

A Mathematical Introduction to Data Science

Yuan Yao

DEPARTMENT OF MATHEMATICS, HONG KONG UNIVERSITY OF SCIENCE AND
TECHNOLOGY, CLEAR WATER BAY, HONG KONG,

E-mail address: `yuany@ust.hk`

URL: https://yao-lab.github.io/book_datasci/

This is a working draft last updated on February 27, 2019.

2000 *Mathematics Subject Classification.* Primary

Key words and phrases. keywords

Special thanks to Amit Singer, Weinan E, Xiuyuan Cheng, and the following students in PKU who help scribe lecture notes with various improvements: Hong Cheng, Chao Deng, Yanzhen Deng, Chendi Huang, Lei Huang, Shujiao Huang, Longlong Jiang, Yuwei Jiang, Wei Jin, Changcheng Li, Xiaoguang Li, Zhen Li, Tengyuan Liang, Feng Lin, Yaning Liu, Peng Luo, Wulin Luo, Tangjie Lv, Yuan Lv, Hongyu Meng, Ping Qin, Jie Ren, Hu Sheng, Zhiming Wang, Yuting Wei, Jiechao Xiong, Jie Xu, Bowei Yan, Jun Yin, and Yue Zhao.

ABSTRACT. This monograph aims to provide graduate students or senior graduates in applied mathematics, computer science and statistics an introduction to data science from a mathematical perspective. The lecture notes have been used in courses at Peking University and Hong Kong University of Science and Technology, that can be found at <https://yao-lab.github.io/course.html>. Some materials are used by Prof. Amit Singer in Princeton University. It is focused on a geometric and topological perspective to data analysis (reduction and visualization), covering a wide range of topics, such as Principal Component Analysis, Multidimensional Scaling, Stein's phenomenon and high dimensionality, shrinkage and regularization, random matrix theory and random projections, robust PCA, sparse PCA, graph realization, manifold learning, and topics in topological data analysis.

Contents

Preface	1
Chapter 1. Geometry of PCA and MDS	3
1. Principal Component Analysis	3
2. How many components? Parallel Analysis	6
3. Multidimensional Scaling	8
4. Theory of MDS (Young/Householder/Schoenberg'1938)	11
5. Reproducing Kernel Hilbert Space and kernel PCA/MDS	13
6. Duality between MDS and PCA as SVD	16
7. Supervised PCA as Sufficient Dimensionality Reduction	17
Exercise	21
Chapter 2. High Dimensional Statistical Models	23
1. Maximum Likelihood Estimate of Mean and Covariance	23
2. Stein's Phenomenon and Shrinkage of Sample Mean	25
3. Random Matrix Theory and Phase Transitions in PCA	35
Chapter 3. Random Projections and Almost Isometry	45
1. Introduction	45
2. The Johnson-Lindenstrauss Lemma	46
3. Example: MDS in Human Genome Diversity Project	49
4. Random Projections and Compressed Sensing	50
Chapter 4. Generalized PCA/MDS via SDP Relaxations	59
1. Introduction of SDP with a Comparison to LP	59
2. Robust PCA	61
3. Exact Recovery Conditions for RPCA	64
4. Tyler's M-estimator	66
5. Sparse PCA	67
6. MDS with Incomplete Information	68
7. Exact Reconstruction and Universal Rigidity	71
8. Maximal Variance Unfolding	73
Chapter 5. Manifold Learning	75
1. Introduction	75
2. ISOMAP	77
3. Locally Linear Embedding (LLE)	80
4. Hessian LLE	84
5. Local Tangent Space Alignment (LTSA)	87
6. Laplacian LLE (Eigenmap)	89
7. Diffusion Map	93

8. Stochastic Neighbor Embedding	95
9. Comparisons	95
Chapter 6. Random Walk on Graphs	97
1. Introduction to Perron-Frobenius Theory and PageRank	97
2. Introduction to Fiedler Theory and Cheeger Inequality	103
3. *Laplacians and the Cheeger inequality for directed graphs	110
4. Lumpability of Markov Chain	116
5. Applications of Lumpability: MNcut and Optimal Reduction of Complex Networks	119
6. Mean First Passage Time	124
7. Transition Path Theory	126
Chapter 7. Diffusion Map	131
1. Diffusion map and Diffusion Distance	131
2. Commute Time Map and Distance	139
3. Diffusion Map: Convergence Theory	142
4. *Vector Diffusion Map	149
Chapter 8. Semi-supervised Learning	157
1. Introduction	157
2. Harmonic Extension of Functions on Graph	157
3. Explanation from Gaussian Markov Random Field	157
4. Explanation from Transition Path Theory	158
5. Well-posedness	159
Chapter 9. Beyond graphs: high dimensional topological/geometric analysis	161
1. From Graph to Simplicial Complex	161
2. Persistent Homology and Discrete Morse Theory	163
3. Exterior Calculus on Complex and Combinatorial Hodge Theory	163
4. Applications of Hodge Theory: Statistical Ranking	165
5. Euler-Calculus	170
Bibliography	173

Preface

This book is used in a course instructed by Yuan Yao at Peking University, part of which is based on a similar course led by Amit Singer at Princeton University.

... the objective of statistical methods is the reduction of data. A quantity of data... is to be replaced by relatively few quantities which shall adequately represent ... the relevant information contained in the original data.

Since the number of independent facts supplied in the data is usually far greater than the number of facts sought, much of the information supplied by an actual sample is irrelevant. It is the object of the statistical process employed in the reduction of data to exclude this irrelevant information, and to isolate the whole of the relevant information contained in the data. —R.A.Fisher

Well, R. A. Fisher exploits a population distribution governed by relatively few “parameters” to characterize the relevant information; in nature relevant vs. irrelevant variations of data can be well described or separated by geometry and even topology, besides statistical models. Especially this is motivated by data visualizations. Therefore, in this course we shall see a dancing among geometric, topological, and statistical data reductions.

CHAPTER 1

Geometry of PCA and MDS

In general, data representation can be vectors, matrices (*esp.* graphs, networks), tensors, and possibly unstructured such as images, videos, languages, sequences, *etc.* In this very first lecture, we start from a basic data representation as Euclidean vectors. Principal Component Analysis (PCA) and Multidimensional Scaling (MDS) are dual one to the other. In PCA, one starts from high dimensional Euclidean representation and looks for a best affine (linear) approximation of data variations; while in MDS, one is exposed with a pairwise distance metric, and pursues an Euclidean representation preserving such a metric. They are actually dual problems, in the sense that given data indeed lying in an Euclidean space, PCA and MDS give different sides of singular vectors of the same data matrix, respectively. Such a geometric picture can be extended to Hilbert spaces of infinite dimension via reproducing kernels as positive definite functions.

1. Principal Component Analysis

Principal component analysis (PCA), invented by Pearson (1901) and Hotelling (1933), is perhaps the most ubiquitous method for dimensionality reduction with high dimensional Euclidean data, under various names in science and engineering such as Karhunen-Lo  e Transform, Empirical Orthogonal Functions, and Principal Orthogonal Decomposition, etc. In the following we will introduce PCA from its geometry.

Let $x_i \in \mathbb{R}^p$, $i = 1, \dots, n$, be n samples in \mathbb{R}^p . Denote the data matrix $X = [x_1 | x_2 | \dots | x_n] \in \mathbb{R}^{p \times n}$. Now we are going to look for a k -dimensional affine space in \mathbb{R}^p to best approximate these n examples (see Figure 1). Assume that such an affine space can be parameterized by $\mu + U\beta$ such that $U = [u_1, \dots, u_k]$ consists of k -columns of an orthonormal basis of the affine space. Then the best approximation in terms of Euclidean distance is given by the following optimization problem.

$$(1) \quad \min_{\beta, \mu, U} I := \sum_{i=1}^n \|x_i - (\mu + U\beta_i)\|^2$$

where $U \in \mathbb{R}^{p \times k}$, $U^T U = I_p$, and $\sum_{i=1}^n \beta_i = 0$ (nonzero sum of β_i can be represented by μ). Taking the first order optimality conditions,

$$\frac{\partial I}{\partial \mu} = -2 \sum_{i=1}^n (x_i - \mu - U\beta_i) = 0 \Rightarrow \hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial I}{\partial \beta_i} = (x_i - \mu - U\beta_i)^T U = 0 \Rightarrow \beta_i = U^T (x_i - \mu)$$

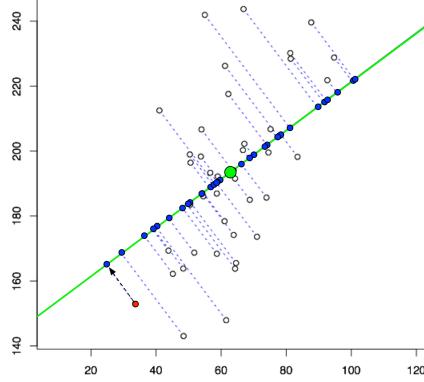


FIGURE 1. Principal Component Analysis as the best affine subspace approximation of data.

Plug in the expression of $\hat{\mu}_n$ and β_i

$$(2) \quad I = \sum_{i=1}^n \|x_i - \hat{\mu}_n - UU^T(x_i - \hat{\mu}_n)\|^2$$

$$(3) \quad = \sum_{i=1}^n \|x_i - \hat{\mu}_n - P_k(x_i - \hat{\mu}_n)\|^2$$

$$(4) \quad = \sum_{i=1}^n \|y_i - P_k(y_i)\|^2, \quad y_i := x_i - \hat{\mu}_n$$

(5)

where $P_k = UU^T$ is a projection operator satisfying the idempotent property $P_k^2 = P_k$.

Denote $Y = [y_1 | y_2 | \cdots | y_n] \in \mathbb{R}^{p \times n}$, whence the original problem turns into

$$\begin{aligned} \min_U \sum_{i=1}^n \|y_i - P_k(y_i)\|^2 &= \min \text{trace}[(Y - P_k Y)^T (Y - P_k Y)] \\ &= \min \text{trace}[Y^T (I - P_k)(I - P_k)Y] \\ &= \min \text{trace}[YY^T(I - P_k)^2] \\ &= \min \text{trace}[YY^T(I - P_k)] \\ &= \min[\text{trace}(YY^T) - \text{trace}(YY^TUU^T)] \\ &= \min[\text{trace}(YY^T) - \text{trace}(U^TYY^TU)]. \end{aligned}$$

Above we use cyclic property of trace and idempotent property of projection.

Since Y does not depend on U , the problem above is equivalent to

$$(6) \quad \max_{UU^T=I_k} \text{Var}(U^T Y) = \max_{UU^T=I_k} \frac{1}{n} \text{trace}(U^T YY^T U) = \max_{UU^T=I_k} \text{trace}(U^T \hat{\Sigma}_n U)$$

where $\hat{\Sigma}_n = \frac{1}{n}YY^T = \frac{1}{n}(X - \hat{\mu}_n\mathbf{1}^T)(X - \hat{\mu}_n\mathbf{1}^T)^T$ is the sample variance¹. Assume that the sample covariance matrix, which is positive semi-definite, has the eigenvalue decomposition $\hat{\Sigma}_n = \hat{U}\hat{\Lambda}\hat{U}^T$, where $\hat{U}^T\hat{U} = I$, $\Lambda = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_n)$, and $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_n \geq 0$. Then

$$\max_{UU^T=I_k} \text{trace}(U^T\hat{\Sigma}_n U) = \sum_{i=1}^k \hat{\lambda}_i$$

In fact when $k = 1$, the maximal covariance is given by the largest eigenvalue along the direction of its associated eigenvector,

$$\max_{\|u\|=1} u^T\hat{\Sigma}_n u =: \hat{\lambda}_1.$$

Restricted on the orthogonal subspace $u \perp \hat{u}_1$ will lead to

$$\max_{\|u\|=1, u^T\hat{u}_1=0} u^T\hat{\Sigma}_n u =: \hat{\lambda}_2,$$

and so on.

Here we conclude that the k -affine space can be discovered by eigenvector decomposition of $\hat{\Sigma}_n$. The sample principal components are defined as column vectors of $\hat{Q} = \hat{U}^TY$, where the j -th observation has its projection on the k -th component as $\hat{q}_k(j) = \hat{u}_k^T y_j = \hat{u}_k^T(x_i - \hat{\mu}_n)$. Therefore, PCA takes the eigenvector decomposition of $\hat{\Sigma}_n = \hat{U}\hat{\Lambda}\hat{U}^T$ and studies the projection of centred data points on top k eigenvectors as the principle components. This is equivalent to the singular value decomposition (SVD) of $X = [x_1, \dots, x_n] \in \mathbb{R}^{p \times n}$ in the following sense,

$$\tilde{X} = XH = X - \frac{1}{n}X \cdot \mathbf{1}\mathbf{1}^T = \tilde{U}\tilde{S}\tilde{V}^T, \quad H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T, \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^n$$

where top *left* singular vectors of centred data matrix $\tilde{X} \in \mathbb{R}^{p \times n}$ are called *principal component loading vectors*, as eigenvectors of the sample covariance matrix $\hat{\Sigma}$. The singular vectors are not unique, but the singular subspace spanned by singular vectors associated with each singular value is unique. Projections of sample point x_i on to the principal component subspace gives the principal component scores, *i.e.* $\beta_i = \tilde{U}_k^T y_i = \tilde{U}_k^T(x_i - \hat{\mu}) \in \mathbb{R}^k$.

How about the *right* singular vectors here? Below we shall see they characterize the metric Multidimensional scaling (MDS). From the properties of singular value decomposition, k -principal components and MDS are thus part of the *best rank- k approximation* of centred data matrix \tilde{X} .

Given a PCA, the following quantities are often used to measure the variances

- total variance:

$$\text{trace}(\hat{\Sigma}_n) = \sum_{i=1}^p \hat{\lambda}_i;$$

¹Note that in statistics the sampled covariance matrix is often defined by for $n \geq 2$,

$$\hat{\Sigma}_n \triangleq \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_n)(x_i - \hat{\mu}_n)^T = \frac{1}{n-1} \tilde{X} \tilde{X}^T.$$

and it becomes 0 when sample size is 1.

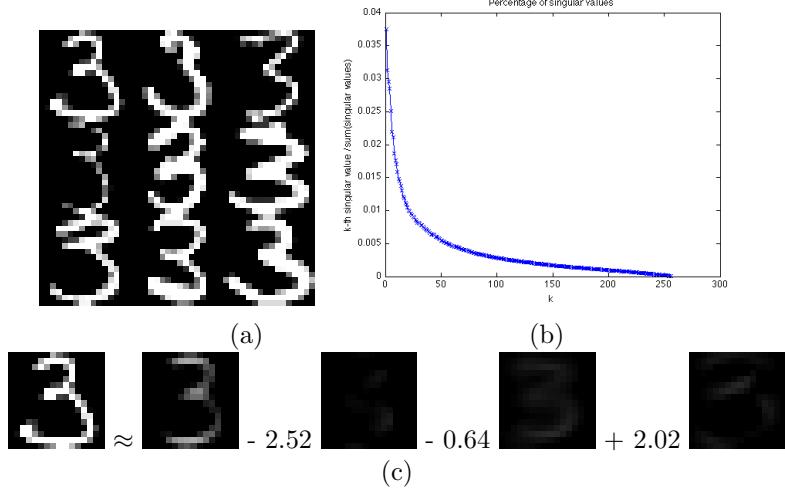


FIGURE 2. (a) random 9 images. (b) percentage of singular values over total sum. (c) approximation of the first image by top 3 principle components (singular vectors).

- percentage of variance explained by top- k principal components:

$$\sum_{i=1}^k \hat{\lambda}_i / \text{trace}(\hat{\Sigma}_n);$$

- generalized variance as total volume:

$$\det(\hat{\Sigma}_n) = \prod_{i=1}^p \hat{\lambda}_i.$$

Example. Take the dataset of hand written digit “3”, $\hat{X} \in \mathbb{R}^{658 \times 256}$ contains 658 images, each of which is of 16-by-16 grayscale image as hand written digit 3. Figure 2 shows a random selection of 9 images, the sorted singular values divided by total sum of singular values, and an approximation of x_1 by top 3 principle components: $x_1 = \hat{\mu}_n - 2.5184\hat{v}_1 - 0.6385\hat{v}_2 + 2.0223\hat{v}_3$.

2. How many components? Parallel Analysis

The key point of PCA is to describe the dataset by only a few principal components (PCs). The question is: how many PCs should be kept? In the example of digit “3”, we have seen that the images of the top eigenvectors are quite relevant, hence are to be kept. The images in Figure 3 are the 201-st to 205-th PCs of the same dataset. They are more like noise, hence are probably not to be kept.

In many occasions, it's sufficient to determine the number of PCs by human experience. An alternative and more precise way is to set a threshold on the percentage of variance. For example, choose k such that

$$\sum_{i=1}^k \hat{\lambda}_i / \text{trace}(\hat{\Sigma}_n) > 0.95.$$

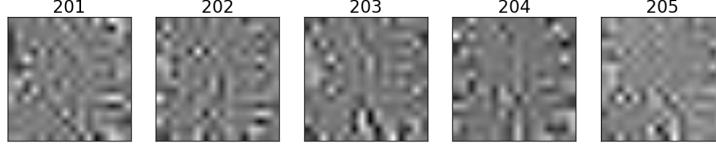


FIGURE 3. The 201-st to 205-th principal components.

Below we shall introduce a new way based on random permutation test, that is called Horn's *parallel analysis* [Hor65, BE92].

2.1. Parallel Analysis for Principal Component Analysis. Take the data matrix $X = [x_1 | x_2 | \cdots | x_n] \in \mathbb{R}^{p \times n}$, generate its parallel data matrices by randomly permuting entries within rows. More precisely, suppose

$$X = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,n} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{p,1} & X_{p,2} & \cdots & X_{p,n} \end{bmatrix}.$$

Randomly take p permutations of n numbers π_1, \dots, π_p (usually π_1 is set as identity). We get a parallel data matrix

$$X^1 = \begin{bmatrix} X_{1,\pi_1(1)} & X_{1,\pi_1(2)} & \cdots & X_{1,\pi_1(n)} \\ X_{2,\pi_2(1)} & X_{2,\pi_2(2)} & \cdots & X_{2,\pi_2(n)} \\ \vdots & \vdots & \ddots & \vdots \\ X_{p,\pi_p(1)} & X_{p,\pi_p(2)} & \cdots & X_{p,\pi_p(n)} \end{bmatrix}.$$

Then we can calculate its singular values $\{\hat{\lambda}_i^1\}_{i=1,\dots,p}$. By choosing a different set of permutations, we can get another parallel data matrix X^2 and its singular values $\{\hat{\lambda}_i^2\}_{i=1,\dots,p}$. Repeat such procedure for R times, we can get R sets of singular values. They can be put together as a matrix

$$\begin{bmatrix} \hat{\lambda}_1^1 & \hat{\lambda}_2^1 & \cdots & \hat{\lambda}_p^1 \\ \hat{\lambda}_1^2 & \hat{\lambda}_2^2 & \cdots & \hat{\lambda}_p^2 \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\lambda}_1^R & \hat{\lambda}_2^R & \cdots & \hat{\lambda}_p^R \end{bmatrix}.$$

For each $i = 1, \dots, p$, compare $\{\hat{\lambda}_i^r\}_{r=1,\dots,R}$ with the i -th singular value $\hat{\lambda}_i$ of the unpermuted data X . Define

$$\text{pval}_i = \frac{1}{R} \#\{\hat{\lambda}_i^r > \hat{\lambda}_i\}.$$

Here $\#$ stands for the cardinality of the set. Notice that $\text{pval}_i \in [0, 1]$. It can be regarded as the probability that $\hat{\lambda}_i$ is indistinguishable from noise. The smaller it is, the more confident we are to think of $\hat{\lambda}_i$ as true signal in the data X . Thus we can set a threshold on $\{\text{pval}_i\}_{i=1,\dots,p}$. For example, we keep $\hat{\lambda}_i$ if $\text{pval}_i < 0.05$.

Let's apply the above parallel analysis to the digit "3" dataset X . The first step is to perform PCA on X , which has been done. The second step is to generate parallel data matrices $\{X^r\}_{r=1,\dots,R}$ ($R = 100$, for example) by randomly permuting

entries within rows of X . Some examples in one of the parallel data matrices are shown in Figure 4. One can see that the randomly permuted images are still informative for digit “3” rather than random images, which implies that each pixel values are highly restricted to some specific domain. This motivates some thoughts that the pixel-wise vectors are restricted on a low-dimensional sub-manifold that will be discussed in later.

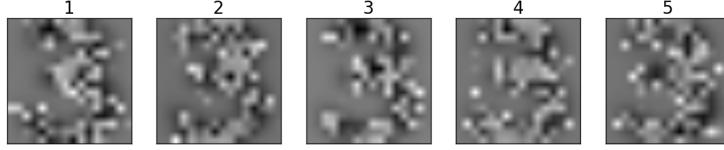


FIGURE 4. Examples of randomly permuted data.

The third step is to calculate the singular values $\{\hat{\lambda}_i^r\}_{i=1,\dots,p}$ of each X^r . The fourth step is to compare them with the corresponding singular values of X , and compute $\{p_{val_i}\}_{i=1,\dots,p}$. Here we choose the threshold $p_{val_i} < 0.05$ for each i . It means that the selected true singular value $\hat{\lambda}_i$ should beat 95% of $\{\hat{\lambda}_i^r\}$. According to the results shown in Figure 5, we can keep the top 19 PCs.

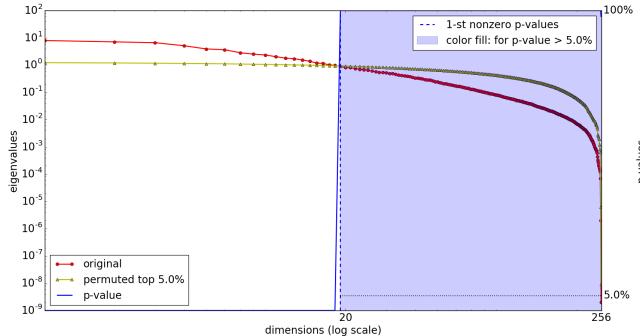


FIGURE 5. Results of parallel analysis on PCA. Considering the exponential decay of eigenvalues and to emphasize the top eigenvalues, log scale are adopted for both axes. The top 5% singular values of the parallel data matrices are draw as reference.

Figure 6 shows the mean image and the top 24 PCs. As we can see, after top 19 PCs, the remaining PCs are still informative for digit “3” as random permuted images are still close to that digit as well. So if the sample points lie on a sub-manifold, permutation test may be conservative in selecting number of PCs.

3. Multidimensional Scaling

Multidimensional Scaling (MDS) roots in psychology [YH41] which aims to recover Euclidean coordinates given pairwise distance metrics or dissimilarities. It is equivalent to PCA when pairwise distances are Euclidean. In the core of

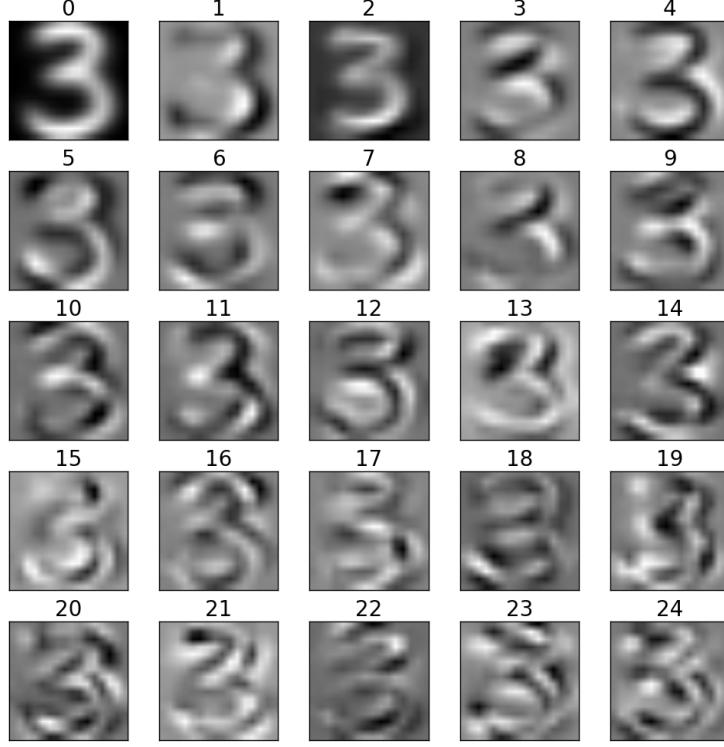


FIGURE 6. Images of the sample mean and the top 24 principal components (top 19 are suggested by parallel analysis). The image No.0 is the sample mean.

theoretical foundation of MDS lies the notion of positive definite functions [Sch37, Sch38a, Sch38b] (or see the survey [Bav11]) which has been the foundation of the kernel method in statistics [Wah90] and modern machine learning society (<http://www.kernel-machines.org/>).

In this section we introduce the classical MDS, or metric embedding problem. The problem of classical MDS or isometric Euclidean embedding is: *given pairwise distances between data points, can one find a system of Euclidean coordinates for those points whose pairwise distances meet given constraints?*

Consider a forward problem: given a set of points $x_1, x_2, \dots, x_n \in \mathbb{R}^p$, let

$$X = [x_1, x_2, \dots, x_n]^{p \times n}.$$

The distance between point x_i and x_j satisfies

$$d_{ij}^2 = \|x_i - x_j\|^2 = (x_i - x_j)^T (x_i - x_j) = x_i^T x_i + x_j^T x_j - 2x_i^T x_j.$$

Now we are considering the inverse problem: *given only d_{ij} , can one find a $\{x_i \in \mathbb{R}^p : i = 1 \dots, n\}$ for some p satisfying the constraint $d_{ij} = \|x_i - x_j\|?$* Clearly the solutions are not unique as any Euclidean transform on $\{x_i\}$ gives another solution. General ideas of classic (metric) MDS is:

- (1) transform squared distance matrix $D = [d_{ij}^2]$ to an inner product form;
- (2) compute the eigen-decomposition for this inner product form.

The key observation is that the two-side centering transform of squared distance matrix D gives the Gram matrix (inner product matrix or kernel matrix) of centred data matrix, i.e.

$$-\frac{1}{2}HDH^T = \tilde{X}^T\tilde{X}.$$

To see this, let K be the inner product or kernel or Gram matrix

$$K = X^TX, \quad X = [x_i] \in \mathbb{R}^{p \times n}$$

with $k = \text{diag}(K_{ii}) \in \mathbb{R}^n$. Note that

$$D = (d_{ij}^2) = k \cdot \mathbf{1}^T + \mathbf{1} \cdot k^T - 2K.$$

where $\mathbf{1} = (1, 1, \dots, 1)^T \in \mathbb{R}^n$.

Now consider the mean and the centred data

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \cdot X \cdot \mathbf{1},$$

$$\tilde{x}_i = x_i - \hat{\mu}_n = x_i - \frac{1}{n} \cdot X \cdot \mathbf{1},$$

or

$$\tilde{X} = X - \frac{1}{n}X \cdot \mathbf{1} \cdot \mathbf{1}^T.$$

In other words, $\tilde{X} = XH$ where $H = I - \frac{1}{n}\mathbf{1} \cdot \mathbf{1}^T$ is the *Householder centering matrix*. Then

$$\tilde{K} \triangleq \tilde{X}^T\tilde{X} = H^TX^TXH = H^TKH.$$

The following lines established the fact that

$$-\frac{1}{2}H \cdot D \cdot H^T = H^TKH.$$

To see this, note that

$$\begin{aligned} B &\triangleq -\frac{1}{2}H \cdot D \cdot H^T \\ &= -\frac{1}{2}H \cdot (k \cdot \mathbf{1}^T + \mathbf{1} \cdot k^T - 2K) \cdot H^T \end{aligned}$$

Since $k \cdot \mathbf{1}^T \cdot H^T = k \cdot \mathbf{1}(I - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^T) = k \cdot \mathbf{1} - k(\frac{\mathbf{1}^T \cdot \mathbf{1}}{n}) \cdot \mathbf{1} = 0$, we have $H \cdot k \cdot \mathbf{1} \cdot H^T = H \cdot \mathbf{1} \cdot k^T \cdot H^T = 0$. This implies that

$$B = -\frac{1}{2}H \cdot D \cdot H^T = H \cdot K \cdot H^T = \tilde{X}^T\tilde{X} = \tilde{K}.$$

Above we have shown that given a squared distance matrix $D = (d_{ij}^2)$, we can convert it to an inner product matrix by $B = -\frac{1}{2}HDH^T$. Eigen-decomposition applied to B will give rise the Euclidean coordinates centred at the origin. Since $B = -\frac{1}{2}HDH^T = \tilde{X}^T\tilde{X} = (XH)^T(XH)$, the eigenvectors of B are exactly right singular vectors of centred data matrix $\tilde{X} = XH$. This shows us the metric MDS is a dual of PCA in the sense that PCA takes left singular vectors of the centred data matrix $\tilde{X} = XH$.

In practice, one often chooses top k nonzero eigenvectors of B for a k -dimensional Euclidean embedding/approximation of data.

Algorithm 1: Classical MDS Algorithm

Input: A squared distance matrix $D^{n \times n}$ with $D_{ij} = d_{ij}^2$.

Output: Euclidean k -dimensional coordinates $\tilde{X}_k \in \mathbb{R}^{k \times n}$ of data.

- 1 Compute $B = -\frac{1}{2}H \cdot D \cdot H^T$, where H is a centering matrix.
- 2 Compute Eigenvalue decomposition $B = U\Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$;
- 3 Choose top k nonzero eigenvalues and corresponding eigenvectors, $\tilde{X}_k = U_k \Lambda_k^{\frac{1}{2}}$ where

$$U_k = [u_1, \dots, u_k], \quad u_k \in \mathbb{R}^n,$$

$$\Lambda_k = \text{diag}(\lambda_1, \dots, \lambda_k)$$

with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$.

In the Classical MDS Algorithm, \tilde{X}_k gives k -dimensional Euclidean coordinations for the n points.

In Matlab, the command for computing classical MDS is "cmdscale", short for Classical Multidimensional Scaling. For non-metric MDS, you may choose "mdscale". Figure 7 shows an example of MDS.

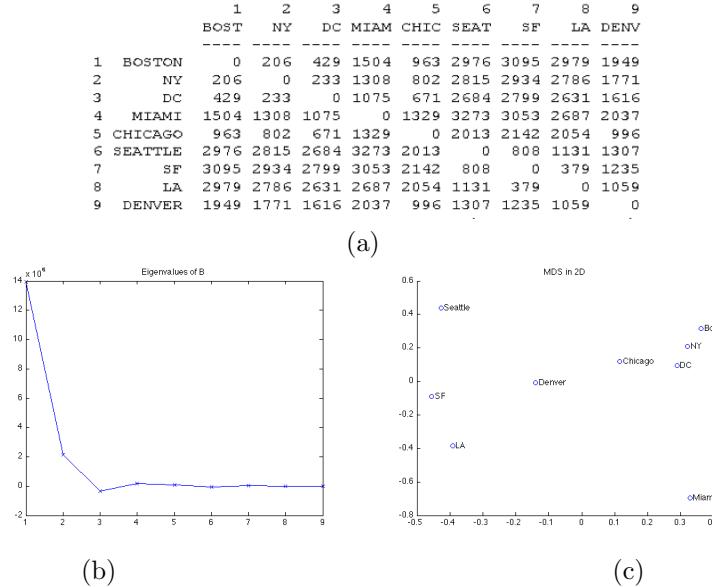


FIGURE 7. MDS of nine cities in USA. (a) Pairwise distances between 9 cities; (b) Eigenvalues of $B = -\frac{1}{2}H \cdot D \cdot H^T$; (c) MDS embedding with top-2 eigenvectors.

4. Theory of MDS (Young/Householder/Schoenberg'1938)

In this section, we shall see that a metric space of n points with distance d_{ij} can be isometrically embedded into an Euclidean space, if and only if the squared distance matrix $D = [d_{ij}^2]$ is conditionally negative definite.

DEFINITION (Positive Semi-definite). Suppose $A^{n \times n}$ is a real symmetric matrix, then A is called positive semi-definite (p.s.d.), denoted by $A \succeq 0$, if $\forall v \in \mathbb{R}^n, v^T A v \geq 0$.

Positive semi-definiteness completely characterizes the inner product matrices: $A \succeq 0 \iff A = Y^T Y$ for some Y .

PROPERTY. Suppose $A^{n \times n}, B^{n \times n}$ are real symmetric matrix, $A \succeq 0, B \succeq 0$. Then we have:

- (1) $A + B \succeq 0$;
- (2) $A \circ B \succeq 0$;

where $A \circ B$ is called Hadamard product and $(A \circ B)_{i,j} := A_{i,j} \times B_{i,j}$.

DEFINITION (Conditionally Negative Definite). Let $A^{n \times n}$ be a real symmetric matrix. A is conditionally negative definite (c.n.d.) $\iff \forall v \in \mathbb{R}^n$, such that $\mathbf{1}^T v = \sum_{i=1}^n v_i = 0$, there holds $v^T A v \leq 0$

The following lemma shows that conditionally negative definite matrices are precisely those matrices whose negative two-side Householder centering transforms are positive semi-definite matrices.

LEMMA 4.1 (Young/Householder-Schoenberg '1938). For any signed probability measure α ($\alpha \in \mathbb{R}^n, \sum_{i=1}^n \alpha_i = 1$),

$$B_\alpha = -\frac{1}{2} H_\alpha C H_\alpha^T \succeq 0 \iff C \text{ is c.n.d.}$$

where H_α is Householder centering matrix: $H_\alpha = \mathbf{I} - \mathbf{1} \cdot \alpha^T$.

PROOF. \Leftarrow We are to show if C is c.n.d., then $B_\alpha \succeq 0$. Taking an arbitrary $x \in \mathbb{R}^n$,

$$x^T B_\alpha x = -\frac{1}{2} x^T H_\alpha C H_\alpha^T x = -\frac{1}{2} (H_\alpha^T x)^T C (H_\alpha^T x).$$

Now we are going to show that $y = H_\alpha^T x$ satisfies $\mathbf{1}^T y = 0$. In fact,

$$\mathbf{1}^T \cdot H_\alpha^T x = \mathbf{1}^T \cdot (\mathbf{I} - \alpha \cdot \mathbf{1}^T) x = (1 - \mathbf{1}^T \cdot \alpha) \mathbf{1}^T \cdot x = 0$$

as $\mathbf{1}^T \cdot \alpha = 1$ for signed probability measure α . Therefore,

$$x^T B_\alpha x = -\frac{1}{2} (H_\alpha^T x)^T C (H_\alpha^T x) \geq 0,$$

as C is c.n.d.

\Rightarrow Now it remains to show if $B_\alpha \succeq 0$ then C is c.n.d. For $\forall x \in \mathbb{R}^n$ satisfying $\mathbf{1}^T \cdot x = 0$, we have

$$H_\alpha^T x = (\mathbf{I} - \alpha \cdot \mathbf{1}^T) x = x - \alpha \cdot \mathbf{1}^T x = x$$

Thus,

$$x^T C x = (H_\alpha^T x)^T C (H_\alpha^T x) = x^T H_\alpha C H_\alpha^T x = -2x^T B_\alpha x \leq 0,$$

as desired.

This completes the proof. □

The following theorem states that conditionally negative definite matrices, after centering, are actually all the squared distance matrices that allows an isometric embedding by the classical MDS algorithm.

THEOREM 4.2 (Classical MDS). Let $D^{n \times n}$ a real symmetric matrix. $C = D - \frac{1}{2}d \cdot \mathbf{1}^T - \frac{1}{2}\mathbf{1} \cdot d^T$, $d = \text{diag}(D)$. Then:

- (1) $B_\alpha = -\frac{1}{2}H_\alpha DH_\alpha^T = -\frac{1}{2}H_\alpha CH_\alpha^T$ for $\forall \alpha$ signed probability measure;
- (2) $C_{i,j} = B_{i,i}(\alpha) + B_{j,j}(\alpha) - 2B_{i,j}(\alpha)$
- (3) D c.n.d. $\iff C$ c.n.d.
- (4) C c.n.d. $\Rightarrow C$ is a square distance matrix (i.e. $\exists Y^{n \times k}$ s.t. $C_{i,j} = \sum_{m=1}^k (y_{i,m} - y_{j,m})^2$)

PROOF. (1) $H_\alpha DH_\alpha^T - H_\alpha CH_\alpha^T = H_\alpha(D - C)H_\alpha^T = H_\alpha(\frac{1}{2}d \cdot \mathbf{1}^T + \frac{1}{2}\mathbf{1} \cdot d^T)H_\alpha^T$.

Since $H_\alpha \cdot \mathbf{1} = 0$, we have

$$H_\alpha DH_\alpha^T - H_\alpha CH_\alpha^T = 0$$

- (2) $B_\alpha = -\frac{1}{2}H_\alpha CH_\alpha^T = -\frac{1}{2}(\mathbf{I} - \mathbf{1} \cdot \alpha^T)C(\mathbf{I} - \alpha \cdot \mathbf{1}^T) = -\frac{1}{2}C + \frac{1}{2}\mathbf{1} \cdot \alpha^T C + \frac{1}{2}C\alpha \cdot \mathbf{1}^T - \frac{1}{2}\mathbf{1} \cdot \alpha^T C\alpha \cdot \mathbf{1}^T$, so we have:

$$B_{i,j}(\alpha) = -\frac{1}{2}C_{i,j} + \frac{1}{2}c_i + \frac{1}{2}c_j - \frac{1}{2}c$$

where $c_i = (\alpha^T C)_i$, $c = \alpha^T C \alpha$. This implies

$$B_{i,i}(\alpha) + B_{j,j}(\alpha) - 2B_{i,j}(\alpha) = -\frac{1}{2}C_{ii} - \frac{1}{2}C_{jj} + C_{ij} = C_{ij},$$

where the last step is due to $C_{i,i} = 0$.

- (3) According to Lemma 4.1 and the first part of Theorem 4.2: C c.n.d. $\iff B$ p.s.d $\iff D$ c.n.d.
- (4) According to Lemma 4.1 and the second part of Theorem 4.2:
 C c.n.d. $\iff B$ p.s.d $\iff \exists Y$ s.t. $B_\alpha = Y^T Y \iff B_{i,j}(\alpha) = \sum_k Y_{i,k} Y_{j,k} \Rightarrow C_{i,j} = \sum_k (Y_{i,k} - Y_{j,k})^2$

This completes the proof. \square

5. Reproducing Kernel Hilbert Space and kernel PCA/MDS

Schoenberg [Sch38b] shows that Euclidean embedding of finite points can be characterized completely by positive definite functions, which paves a way toward Hilbert space embedding. Later Aronzajn [Aro50] developed Reproducing Kernel Hilbert spaces based on positive definite functions which eventually leads to the kernel methods in statistics and machine learning [Vap98, BTA04, CST03].

THEOREM 5.1 (Schoenberg 38). A separable space M with a metric function $d(x, y)$ can be isometrically imbedded in a Hilbert space H , if and only if the family of functions $e^{-\lambda d^2}$ are positive definite for all $\lambda > 0$ (in fact we just need it for a sequence of λ_i whose accumulate point is 0).

Here a symmetric function $k(x, y) = k(y, x)$ is called *positive definite* if for all finite x_i, x_j ,

$$\sum_{i,j} c_i c_j k(x_i, x_j) \geq 0, \quad \forall c_i, c_j$$

with equality holds iff $c_i = c_j = 0$. In other words the function k restricted on $\{(x_i, x_j) : i, j = 1, \dots, n\}$ is a positive definite matrix.

To see the theorem, recall that by the classic MDS theorem, a set of n points with distance d_{ij} can be isometrically imbedded in an Euclidean space if and only if the squared distance matrix is conditionally negative definite, i.e.

$$(7) \quad \sum_{i,j} c_i c_j d_{ij}^2 \leq 0, \quad \sum_i c_i = 0.$$

Using the inverse Fourier transform,

$$(8) \quad e^{-t^2} = \frac{1}{2\sqrt{\pi}} \int_{-\infty}^{\infty} e^{itu} e^{-u^2/4} du$$

function e^{-t^2} is positive definite over the real since it has a positive even spectrum [Sch38b] and so is $e^{-\lambda d^2}$ for $\lambda > 0$ that shows the necessity.

There are two ways to see the sufficiency. First if $e^{-\lambda d^2}$ ($\lambda > 0$) is a positive definite function, then

$$0 \leq \sum_{i,j} c_i c_j \exp(-\lambda d_{ij}^2) = -\lambda \sum_{i,j} c_i c_j d_{ij}^2 + \frac{\lambda^2}{2} \sum_{i,j} c_i c_j d_{ij}^4 - \dots$$

by Tayler expansion. For sufficiently small λ , this implies (7). Since the number of sample points is arbitrary, positive definite function $e^{-\lambda d^2}$ ensures a Hilbert space embedding of possibly infinite dimension.

Second, a powerful observation in [Sch38b] leads to the Schoenberg transform that gives an alternative proof with fruitful results. Notice the formula

$$(9) \quad t^\alpha = c(\alpha) \int_0^\infty (1 - e^{-\lambda^2 t^2}) \lambda^{-1-\alpha} d\lambda,$$

where

$$(10) \quad c(\alpha) = \left[\int_0^\infty (1 - e^{-\lambda^2}) \lambda^{-1-\alpha} \right]^{-1} > 0, \quad 0 < \alpha < 2, t \geq 0.$$

Substituting t by d_{ij} and noticing $\sum_i c_i = 0$, leads to

$$(11) \quad \sum_{i,j} c_i c_j d_{ij}^\alpha = -c(\alpha) \int_0^\infty \sum_{i,j} c_i c_j e^{-\lambda^2 d_{ij}^2} \lambda^{-1-\alpha} d\lambda \leq 0$$

for all $\alpha \in (0, 2)$. Then letting α approach the limit 2 gives the condition (7).

A slightly more general formula than (9) gives the *Schoenberg Transform* Φ as a transform from \mathbb{R}^+ to \mathbb{R}^+ , which takes d to

$$\Phi(d) = \int_0^\infty \frac{1 - \exp(-\lambda d)}{\lambda} g(\lambda) d\lambda,$$

where $g(\lambda)$ is some nonnegative measure on $[0, \infty)$ s.t

$$\int_0^\infty \frac{g(\lambda)}{\lambda} d\lambda < \infty.$$

Sometimes, we may want to transform a square distance matrix to another square distance matrix. The following theorem tells us that Schoenberg Transform [Sch38a, Sch38b] characterizes all the transformations between squared distance matrices.

THEOREM 5.2 (Schoenberg Transform). Given D a squared distance matrix, $C_{i,j} = \Phi(D_{i,j})$. Then

$$C \text{ is a squared distance matrix} \iff \Phi \text{ is a Schoenberg Transform.}$$

Examples of Schoenberg Transforms include

- $\phi_0(d) = d$ with $g_0(\lambda) = \delta(\lambda)$;
- $\phi_1(d) = \frac{1 - \exp(-ad)}{a}$ with $g_1(\lambda) = \delta(\lambda - a)$ ($a > 0$);
- $\phi_2(d) = \ln(1 + d/a)$ with $g_2(\lambda) = \exp(-a\lambda)$;
- $\phi_3(d) = \frac{d}{a(a+d)}$ with $g_3(\lambda) = \lambda \exp(-a\lambda)$;
- $\phi_4(d) = d^p$ ($p \in (0, 1)$) with $g_4(\lambda) = \frac{p}{\Gamma(1-p)} \lambda^{-p}$ (see more in [Bav11]).

The first one gives the identity transform and the last one implies that for a distance function, \sqrt{d} is also a distance function but d^2 is not. To see this, take three points on a line $x = 0, y = 1, z = 2$ where $d(x, y) = d(y, z) = 1$, then for $p > 1$ $d^p(x, z) = 2^p > d^p(x, y) + d^p(y, z) = 2$ which violates the triangle inequality. In fact, d^p ($p \in (0, 1)$) is Euclidean distance function immediately implies the following triangle inequality

$$d^p(0, x+y) \leq d^p(0, x) + d^p(0, y).$$

Note that Schoenberg transform satisfies $\phi(0) = 0$,

$$\begin{aligned} \phi'(d) &= \int_0^\infty \exp(-\lambda d) g(\lambda) d\lambda \geq 0, \\ \phi''(d) &= - \int_0^\infty \exp(-\lambda d) \lambda g(\lambda) d\lambda \leq 0, \end{aligned}$$

and so on. In other words, ϕ is a *completely monotonic function* defined by $(-1)^n \phi^{(n)}(x) \geq 0$, with additional constraint $\phi(0) = 0$. Schoenberg [Sch38a] showed that a function ϕ is completely monotone on $[0, \infty)$ if and only if $\phi(d^2)$ is positive definite and radial on \mathbb{R}^s for all s .

Combined this with Schoenberg transform, one shows that if $d(x, y)$ is an Euclidean distance matrix, then $e^{-\lambda \Phi(d^2)}$ is positive definite for all $\lambda > 0$. Note that for homogeneous function $e^{-\lambda \Phi(tx)} = e^{-\lambda t^k \Phi(x)}$, it suffices to check positive definiteness for $\lambda = 1$.

Symmetric positive definite functions $k(x, y)$ are often called reproducing kernels [Aro50]. In fact the functions spanned by $k_x(\cdot) = k(x, \cdot)$ for $x \in X$ made up of a Hilbert space, where we can associate an inner product induced from $\langle k_x, k_y \rangle = k(x, y)$. The radial basis function $e^{-\lambda d^2} = e^{-\lambda \|x\|^2}$ is often called Gaussian kernel or heat kernel in literature and has been widely used in machine learning.

On the other hand, every Hilbert space \mathcal{H} of functions on \mathcal{X} with bounded evaluation functional can be regarded as a reproducing kernel Hilbert space [Wah90]. By Riesz representation, for every $x \in \mathcal{X}$ there exists $E_x \in \mathcal{H}$ such that $f(x) = \langle f, E_x \rangle$. By boundedness of evaluation functional, $|f(x)| \leq \|f\|_H \|E_x\|$, one can define a reproducing kernel $k(x, y) = \langle E_x, E_y \rangle$ which is bounded, symmetric and positive definite. It is called ‘reproducing’ because we can reproduce the function value using $f(x) = \langle f, k_x \rangle$ where $k_x(\cdot) := k(x, \cdot)$ as a function in \mathcal{H} . Such an universal property makes RKHS a unified tool to study Hilbert function spaces in nonparametric statistics, including Sobolev spaces consisting of splines [Wah90].

Given n samples, the kernel matrix $K = (k(x_i, x_j) : i, j = 1, \dots, n)$ is positive semi-definite matrix, so one can define the following *kernel PCA* which is essentially the *kernel MDS*: *find the top- k eigen-decomposition of the following centred matrix*

$$B = HKH^T$$

then embed the data in the same way as classical MDS.

6. Duality between MDS and PCA as SVD

We have seen that given a set of paired distances d_{ij} , how to find an Euclidean embedding $x_i \in \mathbb{R}^p$ such that $\|x_i - x_j\| = d_{ij}$. However the dimensionality of such an embedding p can be very large. For example, any $n + 1$ points can be isometrically embedded into \mathbb{R}_∞^n using $(d_{i1}, d_{i2}, \dots, d_{in})$ and l_∞ -metric: $d_\infty(x_j, x_k) = \max_{i=1, \dots, n} |d_{ij} - d_{ik}| = d_{ik}$ due to triangle inequality. Moreover, via the heat kernel $e^{-\lambda t^2}$ they can be embedded into Hilbert spaces of infinite dimensions.

Therefore dimensionality reduction is desired when p is large, at the best preservation of pairwise distances.

Given a set of points $x_i \in \mathbb{R}^p$ ($i = 1, 2, \dots, n$); form a data Matrix $X^{p \times n} = [X_1, X_2 \dots X_n]^T$, when p is large, especially in some cases larger than n , we want to find k -dimensional projection with which pairwise distances of the data point are preserved as well as possible. That is to say, if we know the original pairwise distance $d_{ij} = \|X_i - X_j\|$ or data distances with some disturbance $\tilde{d}_{ij} = \|X_i - X_j\| + \epsilon$, we want to find $Y_i \in \mathbb{R}^k$ s.t.:

$$(12) \quad \min_{Y_i \in \mathbb{R}^k} \sum_{i,j} (\|Y_i - Y_j\|^2 - \tilde{d}_{ij}^2)^2.$$

Without loss of generality, we set $\sum_i Y_i = 0$, *i.e.* putting the origin as data center. This is called *nonmetric MDS* for such general \tilde{d}_{ij} .

Schoenberg theory tells us that when d_{ij} is exactly given by Euclidean distance between points, the kernel matrix $B_{ij} = -\frac{1}{2} HDH$ is positive semi-definite, where $D = (d_{ij}^2)$ and $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$, and the optimization problem (12) is equivalent to find $Y_i \in \mathbb{R}^k$ such that

$$(13) \quad \min_{Y \in \mathbb{R}^{k \times n}} \|Y^T Y - B\|_F^2.$$

When B is positive semidefinite, it represents an inner product and therefore the row vectors of matrix Y are the eigenvectors corresponding to k largest eigenvalues of $B = \tilde{X}^T \tilde{X}$, or equivalently the top k *right singular vectors* of $\tilde{X} = USV^T$. This is called *metric MDS*.

We have seen in the first section that the covariance matrix of data $\hat{\Sigma}_n = \frac{1}{n-1} \tilde{X} \tilde{X}^T = \frac{1}{n} US^2 U^T$, passing through the singular vector decomposition (SVD) of $\tilde{X} = USV^T$. Taking top k *left singular vectors* as the embedding coordinates is often called *Principal Component Analysis* (PCA). In PCA, given (centralized) Euclidean coordinate \tilde{X} , usually one gets the inner product matrix as covariance matrix $\hat{\Sigma}_n = \frac{1}{n-1} \tilde{X} \cdot \tilde{X}^T$ which is a $p \times p$ *positive semi-definite* matrix, then the top k eigenvectors of $\hat{\Sigma}_n$ give rise to a k -dimensional embedding of data, as principal components. So both MDS and PCA are unified in SVD of centralized data matrix.

In a summary, consider the data matrix

$$X = [x_1, \dots, x_n]^T \in \mathbb{R}^{n \times p}.^2$$

Let the centred data admits a singular vector decomposition (SVD),

$$\tilde{X} = X - \frac{1}{n}\mathbf{1}\mathbf{1}^T X = HX = \tilde{U}\tilde{S}\tilde{V}^T, \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^n.$$

We have seen that both MDS and PCA can be obtained from such a SVD of centred data matrix.

- MDS embedding is given by top k *left* singular vectors $\tilde{Y}_k^{MDS} = \tilde{U}_k\tilde{S}_k \in \mathbb{R}^{n \times k}$;
- PCA has principle components given by top k *right* singular vectors $\tilde{V}_k \in \mathbb{R}^{p \times k}$ and the projection of centred data on to such a subspace is given by \tilde{Y}_k^{MDS} .

Altogether $\tilde{U}_k\tilde{S}_k\tilde{V}_k^T$ gives best rank- k approximation of \tilde{X} in any unitary invariant norms. In terms of SVD, both PCA and MDS realized the same linear dimensionality reduction in a dual role.

7. Supervised PCA as Sufficient Dimensionality Reduction

Principal component analysis (PCA), invented by Pearson (1901) and Hotelling (1933), is perhaps the most ubiquitous method for dimensionality reduction with high dimensional Euclidean data, under various names in science and engineering such as Karhunen-Loëve Transform, Empirical Orthogonal Functions, and Principal Orthogonal Decomposition, etc. Previously we have seen the geometric interpretation of PCA. It is a type of “unsupervised learning” in the sense that there seems no relation with response variable y . Does PCA really matter with response variable in supervised learning, e.g., in classification or regression?

In the 2005 Fisher Lecture, R. Dennis Cook [Coo07] described PCA as a sufficient dimensionality reduction in regression, and also extended it to principal fitted components (PFC). Here we introduce his idea, together with several variations of supervised PCA: Fisher’s Linear Discriminant Analysis and Li’s Sliced Inverse Regression.

A sufficient dimension reduction Γ ($\Gamma \in \mathbb{R}^{p \times d}$, $\Gamma^T \Gamma = I_d$) refers to the setting that the conditional distribution of $Y|X$ is the same as the distribution of $Y|\Gamma^T X$ for all X .

For example, in regression $Y = f(X, \varepsilon)$, for some unknown function f , sufficient dimensionality reduction implies that $Y = f(\Gamma^T X, \varepsilon)$. However f is unknown here. How can we find Γ independent to the choice of f ?

The answer is a possible Yes when we consider the inverse function, based on the conditional distribution $X|Y$. Consider the scenario where the following inverse model holds: for each value in response variable y ,

$$X_y = \mu + \Gamma \nu_y + \epsilon$$

where for simplicity we assume $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, $\Gamma^T \Gamma = I_d$, and $\sum_y \nu_y = 0$ for removing the degree of freedom in translation. Consider the Inverse Model,

$$X_y = \mu + \Gamma \nu_y + \varepsilon,$$

²In statistics, data matrix is often written as samples row-wise and variables column-wise, i.e. n -by- p matrix. So be careful on your way of writing the data matrix!

where $X_y \in \mathbb{R}^p$, $\nu_y \in \mathbb{R}^d$, $d < p$, the basis $\Gamma \in \mathbb{R}^{p \times d}$ with $\Gamma^T \Gamma = I_d$, and $\varepsilon \sim N_p(0, \sigma^2 I_p)$.

The following proposition states under the assumption of inverse model, Γ is actually a sufficient reduction. See [Coo07] for more details.

PROPOSITION 7.1. Under the inverse model, the distribution of $Y|X$ is the same as the distribution of $Y|\Gamma^T X$.

PROOF. Firstly, $X|Y = y \sim N_p(\mu + \Gamma\nu_y, \sigma^2 I_p)$. By Bayesian formula, we have

$$\begin{aligned} f_{Y|X}(y|x) &\propto f_{X|Y}(x|y)f_Y(y) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}\|x - \mu - \Gamma\nu_y\|^2\right)f_Y(y) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}(\nu_y^T \nu_y - 2\nu_y^T \Gamma^T(x - \mu))\right)f_Y(y) \end{aligned}$$

The last line is given by the orthogonality of Γ . Similarly, since $\Gamma^T X|Y = y \sim N_d(\Gamma^T \mu + \nu_y, \sigma^2 I_d)$, we have

$$\begin{aligned} f_{Y|\Gamma^T X}(y|\Gamma^T x) &\propto f_{\Gamma^T X|Y}(\Gamma^T x|y)f_Y(y) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}\|\Gamma^T x - \Gamma^T \mu - \nu_y\|^2\right)f_Y(y) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}(\nu_y^T \nu_y - 2\nu_y^T \Gamma^T(x - \mu))\right)f_Y(y) \end{aligned}$$

Therefore, the kernel of $Y|X$ and $Y|\Gamma^T X$ are the same, which implies the result. \square

Now consider the Maximum Likelihood Estimate (MLE) of μ , Γ and ν_y . Under the inverse model, the conditional likelihood function

$$f(X_y|\mu, \Gamma, \nu_y) = \frac{1}{\sigma^p \sqrt{(2\pi)^p}} \exp\left[-\frac{1}{2\sigma^2}(X_y - \mu - \Gamma\nu_y)^T(X_y - \mu - \Gamma\nu_y)\right],$$

and the MLE tries to find $\arg \min_{\mu, \Gamma, \nu_y} \prod_y f(X_y|\mu, \Gamma, \nu_y)$, which is equivalent to the following optimization problem after a logarithmic transform:

$$\max_{\mu, \Gamma, \nu_y} -\frac{1}{2\sigma^2} \sum_y \|X_y - \mu - \Gamma\nu_y\|^2 - \sum_y p \log \sigma + C.$$

which leads to the MLE for n examples,

$$\hat{\mu} = \frac{1}{n} \sum_y X_y,$$

$$\nu_y = \hat{\Gamma}^T(X_y - \hat{\mu}),$$

and

$$(14) \quad \hat{\Gamma} = \arg \min_{\Gamma^T \Gamma = I} \sum_y \|X_y - \hat{\mu} - P_\Gamma(X_y - \hat{\mu})\|^2, \quad P_\Gamma = \Gamma \Gamma^T.$$

Comparison with (14) and (2) shows that when y is of distinct values (e.g. the unknown f is injective), this is exactly the PCA in unsupervised learning. Therefore *PCA can be also derived as a sufficient dimensionality reduction in supervised learning, even the function f is unknown here*.

For y with repeated values, Fisher's Linear Discriminant Analysis (LDA) (in classification) [HTF01] and Ker-Chau Li's Sliced Inverse Regression [Li91] (SIR)

are famous candidates for supervised PCAs. Furthermore, Cook [Coo07] a more general class of principal fitted components adapted to supervised learning.

7.1. Linear discriminant analysis. Fisher's *Linear discriminant analysis* (LDA) in classification, like PCA, looks for linear combinations of features which best explain the data. However, LDA attempts to capture the variation between different classes of data.

Algorithm 2: Linear Discriminate Analysis

1 !t

Input: Data with label $\{X_i, y_i\}_{i=1}^N$ where y_i is discrete in $\{1, 2, \dots, K\}$ but not ordered

Output: Effective dimension reducing directions U_d

2 **Step 1:** Compute sample mean and within class means

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i, \quad \hat{\mu}_k = \frac{1}{N_k} \sum_{y_i=k} X_i;$$

Step 2: Compute Between class covariance matrix

$$\hat{\Sigma}_B^{p \times p} = \frac{1}{K} \sum_{k=1}^K (\hat{\mu}_k - \hat{\mu})(\hat{\mu}_k - \hat{\mu})^T;$$

Step 3: Compute Within class covariance matrix

$$\hat{\Sigma}_W^{p \times p} = \frac{1}{N-K} \sum_{k=1}^K \sum_{y_i=k} (X_i - \hat{\mu}_k)(X_i - \hat{\mu}_k)^T;$$

Step 4: Generalized Eigen-decomposition $\hat{\Sigma}_B = \hat{\Sigma}_W U \Lambda U^T$ with

$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$; Choose eigenvectors corresponding to top $d \leq K$ nonzero eigenvalues, i.e., return U_d such that

$$U_d = [u_1, \dots, u_d], \quad u_d \in \mathbb{R}^n.$$

For given data $(X_i, y_i)_{i=1}^N$, assume that $X_i \in \mathbb{R}^p$, and y_i is discrete in $\{1, 2, \dots, K\}$ but not ordered. LDA captures the variance between class and meanwhile discards the variance within class.

Define the *between-class covariance* matrix

$$\hat{\Sigma}_B^{p \times p} = \frac{1}{K} \sum_{k=1}^K (\hat{\mu}_k - \hat{\mu})(\hat{\mu}_k - \hat{\mu})^T;$$

and the *within-class covariance* matrix

$$\hat{\Sigma}_W^{p \times p} = \frac{1}{N-K} \sum_{k=1}^K \sum_{y_i=k} (X_i - \hat{\mu}_k)(X_i - \hat{\mu}_k)^T,$$

where $\hat{\mu}$ is sample mean and $\hat{\mu}_k$ is within class means, i.e.

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{y_i=k} X_i.$$

Now define the Rayleigh quotient by

$$R(w) = \frac{w^T \hat{\Sigma}_B w}{w^T \hat{\Sigma}_W w}$$

which measures, in some sense, the ‘signal-to-noise ratio’ in terms of direction w .

Intuitively, if $\hat{\Sigma}_W$ is invertible, the eigenvector corresponding the largest eigenvalues of $\hat{\Sigma}_W^{-1} \hat{\Sigma}_B$ (or generalized eigenvectors of pair $(\hat{\Sigma}_B, \hat{\Sigma}_W)$) will maximize R . Accordingly, the best feature vectors would be eigenvectors corresponding top k eigenvalues, i.e.

$$U_k = [u_1, u_2, \dots, u_k], \quad u_k \in \mathbb{R}^n,$$

where $\hat{\Sigma}_B u_k = \hat{\Sigma}_W \lambda_k u_k$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$.

NOTE. For *Generalized Eigen Decomposition*(G.E.D) problem $\hat{\Sigma}_B u = \lambda \hat{\Sigma}_W u$, it is more efficient to solve Eigen Decomposition problem $\hat{\Sigma}_W^{-\frac{1}{2}} \hat{\Sigma}_B \hat{\Sigma}_W^{-\frac{1}{2}} \varphi = \lambda \varphi$ first and scale φ by $\hat{\Sigma}_W$, i.e. $u = \hat{\Sigma}_W^{-\frac{1}{2}} \varphi$.

7.2. Sliced Inverse Regression. Ker-Chau Li [Li91] extended LDA from classification to regression by proposing *Sliced Inverse Regression* (SIR). In regression, we are interested in the conditional mean $f(X) = \mathbb{E}[y|X]$, which is a real valued mapping from high dimensional space \mathbb{R}^p and often called the regression function; on the other hand, it is also interesting to look at the inverse conditional mean $g(y) = \mathbb{E}[X|y]$, which is a 1-dimensional curve (manifold) in \mathbb{R}^p often called the principal curve or inverse regression curve. Such a curve might be easier to deal with than the high dimensional regression function.

Algorithm 3: Sliced Inverse Regression

Input: Data with label $\{X_i, y_i\}_{i=1}^N$, where $X_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$ is continuous (or ordered discrete)

Output: Effective dimension reducing directions Γ_d

1 **Step 1:** Divide the range of y_i into S non-overlapping slices $H_s (s = 1, \dots, S)$. N_s is the number of observations within each slice

2 **Step 2:** Compute the sample mean and total covariance matrix

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i, \quad \hat{\Sigma}^{p \times p} = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu})(X_i - \hat{\mu})^T;$$

Step 3: Compute the mean of X_i over all slices and Between slices covariance matrix

$$\hat{\mu}_k = \frac{1}{N_s} \sum_{y_i \in H_s} X_i, \quad \hat{\Sigma}_B^{p \times p} = \frac{1}{K} \sum_h^K (\hat{\mu}_k - \hat{\mu})(\hat{\mu}_k - \hat{\mu})^T;$$

Step 4: Generalized Eigen-decomposition $\hat{\Sigma}_B = \hat{\Sigma} U \Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$; Choose generalized eigenvectors corresponding to top d nonzero eigenvalues, Γ_d i.e.

$$\Gamma_d = [u_1, \dots, u_d], \quad u_k \in \mathbb{R}^n.$$

Given a response variable Y and a random vector $X \in \mathbb{R}^p$ of explanatory variables, SIR is based on the model

$$Y = f(\Gamma X, \epsilon),$$

where $\Gamma^{k \times p}$ is a unknown projection and $k < p$, and f is an unknown link function. One does not need to know f to reconstruct the projection or dimensionality reduction matrix $\Gamma : \mathbb{R}^p \rightarrow \mathbb{R}^d$. In SIR, the range of response values is divided into non-overlapping slices; then one replaces the *between-class covariance* in LDA by the *between-slice covariance*, the *within-class covariance* in LDA by the *total covariance*, respectively; the same generalized eigen-decomposition gives the sufficient dimensionality reduction.

In [WLM09], this algorithm is extended to *Localized Sliced Inverse Regression* (LSIR) which allows for supervised dimension reduction by projection onto a linear subspace that captures the nonlinear subspace relevant to predicting the response.

Algorithm 4: Localized Sliced Inverse Regression

Input: Data with label $\{X_i, y_i\}_{i=1}^N$, where $X_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$ is continuous (or ordered discrete)

Output: Effective dimension reducing directions Γ_d

- 1 **Step 1:** Compute total covariance matrix $\hat{\Sigma}$ as in SIR;
- 2 **Step 2:** Divide the range of y_i into S non-overlapping slices $H_s (s = 1, \dots, S)$; for each sample (X_i, y_i) compute the localized mean

$$\hat{\mu}_{i,loc} = \frac{1}{|s_i|} \sum_{j \in s_i} X_j,$$

where $s_i = \{j : x_j \text{ belongs to the } k \text{ nearest neighbours of } x_i \text{ in } H_s\}$ and s indexes the slice H_s to which i belongs;

- 3 **Step 3:** Compute a localized version of between-slice covariance $\hat{\Sigma}_B$

$$\hat{\Sigma}_{loc} = \frac{1}{N} \sum_i (\hat{\mu}_{i,loc} - \hat{\mu})(\hat{\mu}_{i,loc} - \hat{\mu})^T ;$$

Step 4: Generalized Eigen-decomposition $\hat{\Sigma}_B = \hat{\Sigma} U \Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$; Choose generalized eigenvectors corresponding to top d nonzero eigenvalues, Γ_d i.e.

$$\Gamma_d = [u_1, \dots, u_d], \quad u_k \in \mathbb{R}^n.$$

Exercise

CHAPTER 2

High Dimensional Statistical Models

We have seen that sample mean and covariance in high dimensional Euclidean space \mathbb{R}^p are exploited in Principal Component Analysis (PCA) or its equivalent Multidimensional Scaling (MDS), which are the projections of high dimensional data on to its top singular vectors. In statistics, the sample mean and covariance are the maximal likelihood estimators based on multivariate Gaussian models. In classical statistics with the Law of Large Numbers, for fixed p when sample size $n \rightarrow \infty$, we know such sample mean and covariance will converge, so as to PCA. Although sample mean $\hat{\mu}_n$ and sample covariance $\hat{\Sigma}_n$ are the widely used statistics in multivariate data analysis, they may suffer some problems in high dimensional settings, *e.g.* for large p and small n scenario. In 1956, Stein [Ste56] shows that the sample mean is not the best estimator in terms of *prediction* measured by the mean square error, for $p > 2$; moreover in 2006, Jonestone [Joh06] shows by random matrix theory that PCA might be overwhelmed by random noise for fixed ratio $p/n = \gamma$ when $n \rightarrow \infty$. Among other works, these two pieces of excellent works inspired a long pursuit toward modern high dimensional statistics with a large unexplored field ahead.

1. Maximum Likelihood Estimate of Mean and Covariance

Consider the statistical model $f(X|\theta)$ as a conditional probability function on \mathbb{R}^p with parameter space $\theta \in \Theta$. Let $X_1, \dots, X_n \in \mathbb{R}^p$ are independently and identically distributed (i.i.d.) sampled according to $f(X|\theta_0)$ on \mathbb{R}^p for some $\theta_0 \in \Theta$. The likelihood function is defined as the probability of observing the given data as a function of θ ,

$$L(\theta) = \prod_{i=1}^n f(X_i|\theta),$$

and a maximum likelihood estimator is defined as

$$\hat{\theta}_n^{MLE} \in \arg \max_{\theta \in \Theta} L(\theta) = \arg \max_{\theta \in \Theta} \prod_{i=1}^n f(X_i|\theta)$$

which is equivalent to

$$\arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(X_i|\theta).$$

The following example shows that the sample mean and covariance can be derived from the maximum likelihood estimator under multivariate normal models of data.

1.1. Example: Multivariate Normal Distribution. For example, consider the normal distribution $\mathcal{N}(\mu, \Sigma)$,

$$f(X|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp \left[-\frac{1}{2}(X - \mu)^T \Sigma^{-1} (X - \mu) \right],$$

where $|\Sigma|$ is the determinant of covariance matrix Σ .

To get the MLE of normal distribution, we need to

$$\max_{\mu, \Sigma} P(X_1, \dots, X_n | \mu, \Sigma) = \max_{\mu, \Sigma} \prod_{i=1}^n \frac{1}{\sqrt{2\pi |\Sigma|}} \exp[-(X_i - \mu)^T \Sigma^{-1} (X_i - \mu)]$$

It is equivalent to maximize the log-likelihood

$$I = \log P(X_1, \dots, X_n | \mu, \Sigma) = -\frac{1}{2} \sum_{i=1}^n (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) - \frac{n}{2} \log |\Sigma| + C$$

Let μ^* is the MLE of μ , we have

$$\begin{aligned} 0 &= \frac{\partial I}{\partial \mu^*} = -\sum_{i=1}^n \Sigma^{-1} (X_i - \mu^*) \\ &\Rightarrow \mu^* = \frac{1}{n} \sum_{i=1}^n X_i = \hat{\mu}_n \end{aligned}$$

To get the estimation of Σ , we need to maximize

$$\begin{aligned} I(\Sigma) &= \text{trace}(I) = -\frac{1}{2} \text{trace} \sum_{i=1}^n (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) - \frac{n}{2} \text{trace} \log |\Sigma| + C \\ -\frac{1}{2} \text{trace} \sum_{i=1}^n (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) &= -\frac{1}{2} \sum_{i=1}^n \text{trace}[\Sigma^{-1} (X_i - \mu)(X_i - \mu)^T] \\ &= -\frac{1}{2} (\text{trace} \Sigma^{-1} \hat{\Sigma}_n) (n-1) \\ &= -\frac{n-1}{2} \text{trace}(\Sigma^{-1} \hat{\Sigma}_n^{\frac{1}{2}} \hat{\Sigma}_n^{\frac{1}{2}}) \\ &= -\frac{n-1}{2} \text{trace}(\hat{\Sigma}_n^{\frac{1}{2}} \Sigma^{-1} \hat{\Sigma}_n^{\frac{1}{2}}) \\ &= -\frac{n-1}{2} \text{trace}(S) \end{aligned}$$

where

$$\hat{\Sigma}_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}_n)(X_i - \hat{\mu}_n)^T,$$

$S = \hat{\Sigma}_n^{\frac{1}{2}} \Sigma^{-1} \hat{\Sigma}_n^{\frac{1}{2}}$ is symmetric and positive definite. Above we repeatedly use cyclic property of trace:

- $\text{trace}(AB) = \text{trace}(BA)$, or more generally
- (invariance under cyclic permutation group) $\text{trace}(ABCD) = \text{trace}(BCDA) = \text{trace}(CDBA) = \text{trace}(DABC)$.

Then we have

$$\Sigma = \hat{\Sigma}_n^{-\frac{1}{2}} S^{-1} \hat{\Sigma}_n^{-\frac{1}{2}}$$

$$-\frac{n}{2} \log |\Sigma| = \frac{n}{2} \log |S| + \frac{n}{2} \log |\hat{\Sigma}_n| = f(\hat{\Sigma}_n)$$

Therefore,

$$\max I(\Sigma) \Leftrightarrow \min \frac{n-1}{2} \text{trace}(S) - \frac{n}{2} \log |S| + \text{Const}(\hat{\Sigma}_n, 1)$$

Suppose $S = U\Lambda U$ is the eigenvalue decomposition of S , $\Lambda = \text{diag}(\lambda_i)$

$$J = \frac{n-1}{2} \sum_{i=1}^p \lambda_i - \frac{n}{2} \sum_{i=1}^p \log(\lambda_i) + \text{Const}$$

$$\frac{\partial J}{\partial \lambda_i} = \frac{n-1}{2} - \frac{n}{2} \frac{1}{\lambda_i} \Rightarrow \lambda_i = \frac{n}{n-1}$$

$$S = \frac{n}{n-1} I_p$$

This gives the MLE solution

$$\Sigma^* = \frac{n-1}{n} \hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)(X_i - \hat{\mu}_n)^T,$$

which differs to $\hat{\Sigma}_n$ only in that the denominator $(n-1)$ is replaced by n . In covariance matrix, $(n-1)$ is used because for a single sample $n=1$, there is no variance at all.

Under some regularity conditions, the maximum likelihood estimator $\hat{\theta}_n^{MLE}$ has the following nice *limiting* properties:

- A. (Consistency) $\hat{\theta}_n^{MLE} \rightarrow \theta_0$, in probability and almost surely.
- B. (Asymptotic Normality) $\sqrt{n}(\hat{\theta}_n^{MLE} - \theta_0) \rightarrow \mathcal{N}(0, I_0^{-1})$ in distribution, where I_0 is the Fisher Information matrix

$$I(\theta_0) := \mathbb{E}\left[\left(\frac{\partial}{\partial \theta} \log f(X|\theta_0)\right)^2\right] = -\mathbb{E}\left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta_0)\right].$$

- C. (Asymptotic Efficiency) $\lim_{n \rightarrow \infty} \text{cov}(\hat{\theta}_n^{MLE}) = I^{-1}(\theta_0)$. Hence $\hat{\theta}_n^{MLE}$ is the Uniformly Minimum-Variance Unbiased Estimator, i.e. the estimator with the least variance among the class of unbiased estimators, for any unbiased estimator $\hat{\theta}_n$, $\lim_{n \rightarrow \infty} \text{var}(\hat{\theta}_n^{MLE}) \leq \lim_{n \rightarrow \infty} \text{var}(\hat{\theta}_n)$.

The asymptotic results all hold under the assumption by fixing p and taking $n \rightarrow \infty$, where MLE satisfies $\hat{\mu}_n \rightarrow \mu$ and $\hat{\Sigma}_n \rightarrow \Sigma$. However as we can see in the following, $\hat{\mu}_n$ is not the best estimators for *prediction* when the dimension of the data p gets large, with a finite sample n .

2. Stein's Phenomenon and Shrinkage of Sample Mean

2.1. Risk and Bias-Variance Decomposition. To measure the *prediction* performance of an estimator $\hat{\mu}_n$, it is natural to consider the expected squared loss in regression, i.e. given a response $y = \mu + \epsilon$ with zero mean noise $\mathbb{E}\epsilon = 0$,

$$\mathbb{E}\|y - \hat{\mu}_n\|^2 = \mathbb{E}\|\mu - \hat{\mu} + \epsilon\|^2 = \mathbb{E}\|\mu - \hat{\mu}\|^2 + \text{Var}(\epsilon), \quad \text{Var}(\epsilon) = \mathbb{E}(\epsilon^T \epsilon).$$

Since $Var(\epsilon)$ is a constant for all estimators $\hat{\mu}$, one may simply look at the first part which is often called as *risk* in literature,

$$R(\hat{\mu}_n, \mu) = \mathbb{E}L(\hat{\mu}_n, \mu)$$

where the loss function takes the square loss here,

$$L(\hat{\mu}_n, \mu) = \|\hat{\mu}_n - \mu\|^2.$$

It is the *mean square error* (MSE) between μ and its estimator $\hat{\mu}$.

The risk or MSE enjoy the following important *bias-variance decomposition*, as a result of the Pythagorean theorem.

$$\begin{aligned} R(\hat{\mu}_n, \mu) &= \mathbb{E}\|\hat{\mu}_n - \mathbb{E}[\hat{\mu}_n] + \mathbb{E}[\hat{\mu}_n] - \mu\|^2 \\ &= \mathbb{E}\|\hat{\mu}_n - \mathbb{E}[\hat{\mu}_n]\|^2 + \|\mathbb{E}[\hat{\mu}_n] - \mu\|^2 \\ &=: Var(\hat{\mu}_n) + Bias(\hat{\mu}_n)^2 \end{aligned}$$

Consider multivariate Gaussian model: let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \Sigma)$, $X_i \in \mathbb{R}^p (i = 1 \dots n)$, then the maximum likelihood estimators (MLE) of the parameters (μ and Σ) are as follows:

$$\hat{\mu}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \hat{\Sigma}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)(X_i - \hat{\mu}_n)^T.$$

For simplicity, take a coordinate transform (PCA) $Y_i = U^T X_i$ where $\Sigma = U \Lambda U^T$ is an eigen-decomposition of the population covariance matrix Σ . Assume that $\Lambda = \sigma^2 I_p$ and $n = 1$, then it suffices to consider $Y \sim \mathcal{N}(\mu, \sigma^2 I_p)$ in the sequel. In this case $\hat{\mu}_n^{MLE} = Y$.

The following example shows the bias and variance of MLE.

EXAMPLE 1. For the simple case $Y_i \sim \mathcal{N}(\mu, \sigma^2 I_p)$ ($i = 1, \dots, n$), the MLE estimator satisfies

$$Bias(\hat{\mu}_n^{MLE}) = 0$$

and

$$Var(\hat{\mu}_n^{MLE}) = \frac{p}{n} \sigma^2$$

In particular for $n = 1$, $Var(\hat{\mu}_n^{MLE}) = \sigma^2 p$ for $\hat{\mu}_n^{MLE} = Y$.

EXAMPLE 2. MSE of Linear Estimators. Consider $Y \sim \mathcal{N}(\mu, \sigma^2 I_p)$ and linear estimator $\hat{\mu}_C = CY$. Then we have

$$Bias(\hat{\mu}_C) = \|(I - C)\mu\|^2$$

and

$$Var(\hat{\mu}_C) = \mathbb{E}[(CY - C\mu)^T(CY - C\mu)] = \mathbb{E}[\text{trace}((Y - \mu)^T C^T C(Y - \mu))] = \sigma^2 \text{trace}(C^T C).$$

Linear estimator includes an important case, the *Ridge regression* (also known as Tikhonov regularization in applied mathematics) with $C = X(X^T X + \lambda I)^{-1} X^T$,

$$\min_{\mu} \frac{1}{2} \|Y - X\beta\|^2 + \frac{\lambda}{2} \|\beta\|^2, \quad \lambda > 0.$$

For simplicity, one may restrict our discussions on the diagonal linear estimators $C = \text{diag}(c_i)$ (up to an change of orthonormal basis for Ridge regression), whose risk is

$$R(\hat{\mu}_C, \mu) = \sigma^2 \sum_{i=1}^p c_i^2 + \sum_{i=1}^p (1 - c_i)^2 \mu_i^2.$$

In this case, it is simple to find minimax risk over the hyper-rectangular model class $|\mu_i| \leq \tau_i$,

$$\inf_{c_i} \sup_{|\mu_i| \leq \tau_i} R(\hat{\mu}_C, \mu) = \sum_{i=1}^p \frac{\sigma^2 \tau_i^2}{\sigma^2 + \tau_i^2}.$$

From here one can see that for those sparse model classes such that $\#\{i : \tau_i = O(\sigma)\} = k \ll p$, it is possible to get smaller risk using linear estimators than MLE.

In general, is it possible to introduce some *biased* estimators which significantly reduces the *variance* such that the total risk is smaller than MLE uniformly for all μ ? This is the notion of inadmissibility introduced by Charles Stein in 1956 and he find the answer is YES by presenting the James-Stein estimators, as the shrinkage of sample means.

DEFINITION (Inadmissible). An estimator $\hat{\mu}_n$ of the parameter μ is called **inadmissible** on \mathbb{R}^p with respect to the squared risk if there exists another estimator μ_n^* such that

$$\mathbb{E}\|\mu_n^* - \mu\|^2 \leq \mathbb{E}\|\hat{\mu}_n - \mu\|^2 \quad \text{for all } \mu \in \mathbb{R}^p,$$

and there exist $\mu_0 \in \mathbb{R}^p$ such that

$$\mathbb{E}\|\mu_n^* - \mu_0\|^2 < \mathbb{E}\|\hat{\mu}_n - \mu_0\|^2.$$

In this case, we also call that μ_n^* **dominates** $\hat{\mu}_n$. Otherwise, the estimator $\hat{\mu}_n$ is called **admissible**.

The notion of inadmissibility or dominance introduces a partial order on the set of estimators where admissible estimators are local optima in this partial order.

Stein (1956) [Ste56] found that if $p \geq 3$, then the MLE estimator $\hat{\mu}_n$ is inadmissible. This property is known as **Stein's phenomenon**. This phenomenon can be described like:

For $p \geq 3$, there exists $\hat{\mu}$ such that $\forall \mu$,

$$R(\hat{\mu}, \mu) < R(\hat{\mu}^{MLE}, \mu)$$

which makes MLE inadmissible.

A typical choice is the *James-Stein estimator*.

EXAMPLE 3 (James-Stein Estimator). Charles Stein shows in 1956 that MLE is inadmissible, while the following original form of James-Stein estimator is demonstrated by his student Willard James in 1961. Bradley Efron [Efr10] summarizes the history and gives a simple derivation of these estimators from an Empirical Bayes point of view.

$$(15) \quad \hat{\mu}^{JS} = \left(1 - \frac{\sigma^2(p-2)}{\|\hat{\mu}^{MLE}\|}\right) \hat{\mu}^{MLE}.$$

Such an estimator shrinks each component of $\hat{\mu}^{MLE}$ toward 0. However, one can shrink it toward other point such as the mean component of $\hat{\mu}^{MLE}$: $\bar{z} = \sum_{i=1}^p z_i/p$ and $S(z) := \sum_i (z_i - \bar{z})^2$,

$$(16) \quad \hat{\mu}_i^{JS+} = \bar{z} + \left(1 - \frac{\sigma^2(p-3)}{S(\hat{\mu}^{MLE})}\right) \hat{\mu}_i^{MLE},$$

at the sacrifice of $p \geq 4$ to reach the same phenomenon.

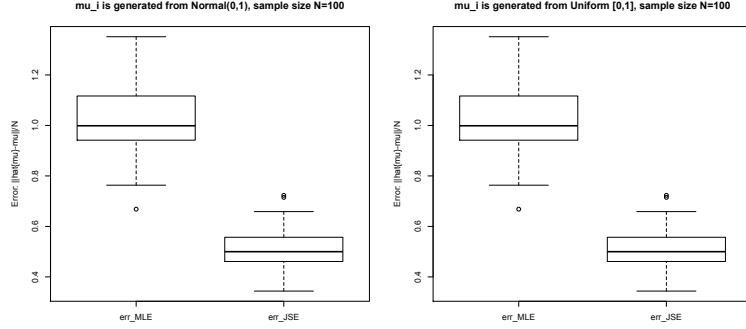


FIGURE 1. Comparison of risks between Maximum Likelihood Estimators and James-Stein Estimators.

THEOREM 2.1. Suppose $Y \sim \mathcal{N}_p(\mu, I)$. Then $\hat{\mu}^{\text{MLE}} = Y$. $R(\hat{\mu}, \mu) = \mathbb{E}_{\mu} \|\hat{\mu} - \mu\|^2$, and define

$$\hat{\mu}^{JS} = \left(1 - \frac{p-2}{\|Y\|^2}\right) Y$$

then

$$R(\hat{\mu}^{JS}, \mu) < R(\hat{\mu}^{\text{MLE}}, \mu)$$

Figure 1 shows some simulations where James-Stein Estimator dominates Maximum Likelihood Estimator. Table 1 gives a real world example by Bradley Efron showing that JSEs improve MLE, yet for some extreme case like Clemente, JSE may suffer from over-shrinkage toward the average.

Next we outline the proof of such results. First of all, we'll prove a useful lemma.

2.2. Stein's Unbiased Risk Estimates (SURE). Discussions below are all under the assumption that $Y \sim \mathcal{N}_p(\mu, I)$.

LEMMA 2.2. (Stein's Unbiased Risk Estimates (SURE)) Suppose $\hat{\mu} = Y + g(Y)$, g satisfies ¹

- (1) g is weakly differentiable.
- (2) $\sum_{i=1}^p \int |\partial_i g_i(x)| dx < \infty$

then

$$(17) \quad R(\hat{\mu}, \mu) = \mathbb{E}_{\mu} (p + 2\nabla^T g(Y) + \|g(Y)\|^2)$$

where $\nabla^T g(Y) := \sum_{i=1}^p \frac{\partial}{\partial y_i} g_i(Y)$.

Examples of $g(x)$: For James-Stein estimator

$$g(x) = -\frac{p-2}{\|Y\|^2} Y$$

and for soft-thresholding, each component

$$g_i(x) = \begin{cases} -\lambda & x_i > \lambda \\ -x_i & |x_i| \leq \lambda \\ \lambda & x_i < -\lambda \end{cases}$$

¹cf. p38, Prop 2.4 [GE]

TABLE 1. Efron's Battting example. There are $p = 18$ players and $n = 45$ samples in the early part of 1970 season. $\hat{\mu}^{MLE}$ is obtained from the mean hits in these early games, while μ is obtained by averages over the remainder of the season. $\hat{\mu}^{JS0}$ takes the shrinkage toward 0 (Eq.(15)), and $\hat{\mu}^{JS+}$ takes the shrinkage toward the average $\bar{\mu}^{MLE}$ (Eq. (16)). Both forms of JS-estimators improve MLE by reducing the mean square error in prediction, while the latter enjoys a much more noticeable improvement than the former.

Name	hits/AB	$\hat{\mu}_i^{(MLE)}$	μ_i	$\hat{\mu}_i^{(JS0)}$	$\hat{\mu}_i^{(JS+)}$
Clemente	18/45	0.4	0.346	0.378	0.294
F.Robinson	17/45	0.378	0.298	0.357	0.289
F.Howard	16/45	0.356	0.276	0.336	0.285
Johnstone	15/45	0.333	0.222	0.315	0.28
Berry	14/45	0.311	0.273	0.294	0.275
Spencer	14/45	0.311	0.27	0.294	0.275
Kessinger	13/45	0.289	0.263	0.273	0.27
L.Alvarado	12/45	0.267	0.21	0.252	0.266
Santo	11/45	0.244	0.269	0.231	0.261
Swoboda	11/45	0.244	0.23	0.231	0.261
Unser	10/45	0.222	0.264	0.21	0.256
Williams	10/45	0.222	0.256	0.21	0.256
Scott	10/45	0.222	0.303	0.21	0.256
Petrocelli	10/45	0.222	0.264	0.21	0.256
E.Rodriguez	10/45	0.222	0.226	0.21	0.256
Campaneris	9/45	0.2	0.286	0.189	0.252
Munson	8/45	0.178	0.316	0.168	0.247
Alvis	7/45	0.156	0.2	0.147	0.242
Mean Square Error	-	0.075545	-	0.072055	0.021387

Both of them are weakly differentiable. But Hard-Thresholding:

$$g_i(x) = \begin{cases} 0 & |x_i| > \lambda \\ -x_i & |x_i| \leq \lambda \end{cases}$$

which is not weakly differentiable!

PROOF. Let $\phi(y)$ be the density function of standard Normal distribution $\mathcal{N}_p(0, I)$.

$$\begin{aligned} R(\hat{\mu}, \mu) &= \mathbb{E}_{\mu} \|Y + g(Y) - \mu\|^2 \\ &= \mathbb{E}_{\mu} (p + 2(Y - \mu)^T g(Y) + \|g(Y)\|^2) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mu}(Y - \mu)^T g(Y) &= \sum_{i=1}^p \int_{-\infty}^{\infty} (y_i - \mu_i) g_i(Y) \phi(Y - \mu) dY \\ &= \sum_{i=1}^p \int_{-\infty}^{\infty} -g_i(Y) \frac{\partial}{\partial y_i} \phi(Y - \mu) dY, \text{ derivative of Gaussian function} \\ &= \sum_{i=1}^p \int_{-\infty}^{\infty} \frac{\partial}{\partial y_i} g_i(Y) \phi(Y - \mu) dY, \text{ Integration by parts} \\ &= \mathbb{E}_{\mu} \nabla^T g(Y) \end{aligned}$$

□

Thus, we define

$$(18) \quad U(Y) := p + 2\nabla^T g(Y) + \|g(Y)\|^2$$

for convenience, and $R(\hat{\mu}, \mu) = \mathbb{E}_\mu U(Y)$.

This lemma is in fact called the Stein's lemma in Tsybakov's book [Tsy09] (page 157~158).

2.3. Risk of Linear Estimator.

$$\begin{aligned} \hat{\mu}_C(Y) &= Cy \\ g(Y) &= (C - I)Y \\ \nabla^T g(Y) &= -\sum_i \frac{\partial}{\partial y_i} ((C - I)Y) = \text{trace}(C) - p \\ U(Y) &= p + 2\nabla^T g(Y) + \|g(Y)\|^2 \\ &= p + 2(\text{trace}(C) - p) + \|(I - C)Y\|^2 \\ &= -p + 2\text{trace}(C) + \|(I - C)Y\|^2 \end{aligned}$$

In applications, $C = C(\lambda)$ often depends on some regularization parameter λ (e.g. ridge regression). So one could find optimal λ^* by minimizing the MSE over λ . Suppose $Y \sim \mathcal{N}(\mu, \sigma^2 I)$,

$$R(\hat{\mu}_C, \mu) = \|(I - C(\lambda))Y\|^2 - p\sigma^2 + 2\sigma^2\text{trace}(C(\lambda)).$$

2.4. Risk of James-Stein Estimator. Recall

$$\begin{aligned} g(Y) &= -\frac{p-2}{\|Y\|^2}Y \\ U(Y) &= p + 2\nabla^T g(Y) + \|g(Y)\|^2 \\ \|g(Y)\|^2 &= \frac{(p-2)^2}{\|Y\|^2} \\ \nabla^T g(Y) &= -\sum_i \frac{\partial}{\partial y_i} \left(\frac{p-2}{\|Y\|^2} Y \right) = -\frac{(p-2)^2}{\|Y\|^2} \end{aligned}$$

we have

$$R(\hat{\mu}^{\text{JS}}, \mu) = \mathbb{E}U(Y) = p - \mathbb{E}_\mu \frac{(p-2)^2}{\|Y\|^2} < p = R(\hat{\mu}^{\text{MLE}}, \mu)$$

when $p \geq 3$.

PROBLEM. What's wrong when $p = 1$? Does SURE still hold?

REMARK. Indeed, we have the following theorem

THEOREM 2.3 (Lemma 2.8 in Johnstone's book (GE)). $Y \sim N(\mu, I)$, $\forall \hat{\mu} = CY$, $\hat{\mu}$ is admissible iff

- (1) C is symmetric.
- (2) $0 \leq \rho_i(C) \leq 1$ (eigenvalue).
- (3) $\rho_i(C) = 1$ for at most two i .

To find an upper bound of the risk of James-Stein estimator, notice that $\|Y\|^2 \sim \chi^2(\|\mu\|^2, p)$ and ²

$$\chi^2(\|\mu\|^2, p) \stackrel{d}{=} \chi^2(0, p + 2N), \quad N \sim \text{Poisson} \left(\frac{\|\mu\|^2}{2} \right)$$

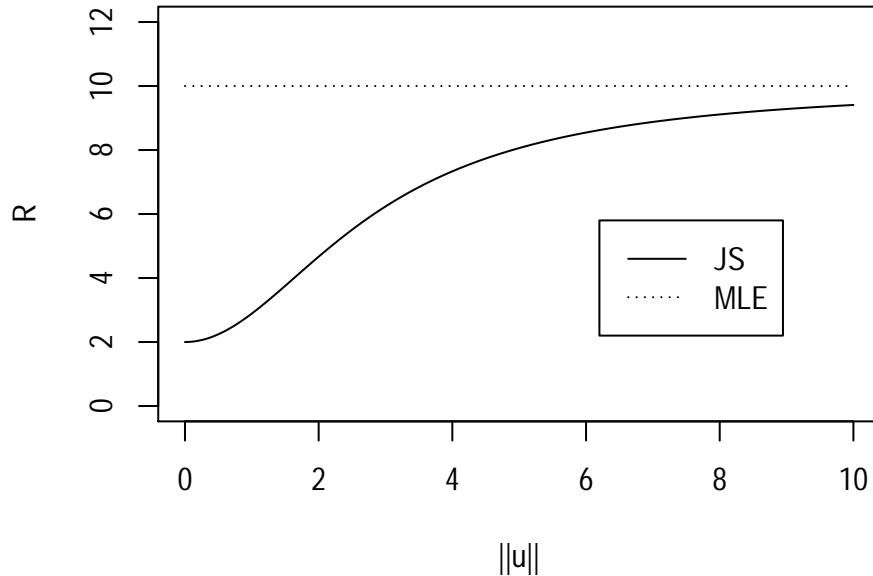
we have

$$\begin{aligned} \mathbb{E}_\mu \left(\frac{1}{\|Y\|^2} \right) &= \mathbb{E} \mathbb{E}_\mu \left[\frac{1}{\|Y\|^2} \mid N \right] \\ &= \mathbb{E} \frac{1}{p + 2N - 2} \\ &\geq \frac{1}{p + 2\mathbb{E}N - 2} \quad (\text{Jensen's Inequality}) \\ &= \frac{1}{p + \|\mu\|^2 - 2} \end{aligned}$$

that is

PROPOSITION 2.4 (Upper bound of MSE for the James-Stein Estimator). $Y \sim \mathcal{N}(\mu, I_p)$,

$$R(\hat{\mu}^{\text{JS}}, \mu) \leq p - \frac{(p-2)^2}{p-2+\|\mu\|^2} = 2 + \frac{(p-2)\|\mu\|^2}{p-2+\|\mu\|^2}$$



²This is a homework.

2.5. Risk of Soft-thresholding. Using Stein's unbiased risk estimate, we have soft-thresholding in the form of

$$\hat{\mu}(x) = x + g(x). \quad \frac{\partial}{\partial i} g_i(x) = -I(|x_i| \leq \lambda)$$

We then have

$$\begin{aligned} \mathbb{E}_\mu \|\hat{\mu}_\lambda - \mu\|^2 &= \mathbb{E}_\mu \left(p - 2 \sum_{i=1}^p I(|x_i| \leq \lambda) + \sum_{i=1}^p x_i^2 \wedge \lambda^2 \right) \\ &\leq 1 + (2 \log p + 1) \sum_{i=1}^p \mu_i^2 \wedge 1 \quad \text{if we take } \lambda = \sqrt{2 \log p} \end{aligned}$$

By using the inequality

$$\frac{1}{2}a \wedge b \leq \frac{ab}{a+b} \leq a \wedge b$$

we can compare the risk of soft-thresholding and James-Stein estimator as

$$1 + (2 \log p + 1) \sum_{i=1}^p (\mu_i^2 \wedge 1) \leq 2 + c \left(\left(\sum_{i=1}^p \mu_i^2 \right) \wedge p \right) \quad c \in (1/2, 1)$$

In LHS, the risk for each μ_i is bounded by 1 so if μ is sparse ($s = \#\{i : \mu_i \neq 0\}$) but large in magnitudes (s.t. $\|\mu\|_2^2 \geq p$), we may expect LHS = $O(s \log p) < O(p) =$ RHS.³

In addition to L_1 penalty in LASSO, there are also other penalty functions like

- $\lambda \|\beta\|_0$ This leads to *hard-thresholding* when $X = I$. Solving this problem is normally NP-hard.
 - $\lambda \|\beta\|_p$, $0 < p < 1$. Non-convex, also NP-hard.
 - $\lambda \sum \rho(\beta_i)$. such that
 - (1) $\rho'(0)$ singular (for sparsity in variable selection)
 - (2) $\rho'(\infty) = 0$ (for unbiasedness in parameter estimation)
- Such ρ must be non-convex essentially (Jianqing Fan and Runze Li, 2001).

2.6. Appendix: more details on deriving James-Stein estimator. Now, let us look for a function g such that the risk of the estimator $\tilde{\mu}_n(Y) = (1 - g(Y))Y$ is smaller than the MLE of $Y \sim \mathcal{N}(\mu, \sigma^2 I_p)$. We have

$$\begin{aligned} \mathbb{E}\|\tilde{\mu}_n - \mu\|^2 &= \sum_{i=1}^p \mathbb{E}[(1 - g(y))y_i - \mu_i]^2 \\ &= \sum_{i=1}^p \{\mathbb{E}[(y_i - \mu_i)^2] + 2\mathbb{E}[(\mu_i - y_i)g(y)y_i] \\ &\quad + \mathbb{E}[y_i^2 g(y)^2]\}. \end{aligned}$$

Suppose now that the function g is such that the assumptions of Stein's Lemma 2.5 hold (page 157~158 in Tsybakov's book [Tsy09]), i.e. weakly differentiable.

LEMMA 2.5 (Stein's lemma). Suppose that a function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ satisfies:

- (i) $f(u_1, \dots, u_p)$ is absolutely continuous in each coordinate u_i for almost all values (with respect to the Lebesgue measure on \mathbb{R}^{p-1}) of other coordinates $(u_j, j \neq i)$

³also cf. p43 [GE]

(ii)

$$\mathbb{E}\left|\frac{\partial f(y)}{\partial y_i}\right| < \infty, \quad i = 1, \dots, p.$$

then

$$\mathbb{E}[(\mu_i - y_i)f(y)] = -\varepsilon^2 \mathbb{E}\left[\frac{\partial f}{\partial y_i}(y)\right], \quad i = 1, \dots, p.$$

With Stein's Lemma, therefore

$$\mathbb{E}[(\mu_i - y_i)(1 - g(y))y_i] = -\varepsilon^2 \mathbb{E}\left[g(y) + y_i \frac{\partial g}{\partial y_i}(y)\right],$$

with

$$\mathbb{E}[(y_i - \mu_i)^2] = \varepsilon^2 = \sigma^2,$$

we have

$$\mathbb{E}[(\tilde{\mu}_{n,i} - \mu_i)^2] = \varepsilon^2 - 2\varepsilon^2 \mathbb{E}\left[g(y) + y_i \frac{\partial g}{\partial y_i}(y)\right] + \mathbb{E}[y_i^2 g(y)^2].$$

Summing over i gives

$$\mathbb{E}\|\tilde{\mu}_n - \mu\|^2 = p\varepsilon^2 + \mathbb{E}[W(y)] = \mathbb{E}\|\hat{\mu}_n - \mu\|^2 + \mathbb{E}[W(y)]$$

with

$$W(y) = -2p\varepsilon^2 g(y) + 2\varepsilon^2 \sum_{i=1}^p y_i \frac{\partial g}{\partial y_i}(y) + \|y\|^2 g(y)^2.$$

The risk of $\tilde{\mu}_n$ is smaller than that of $\hat{\mu}_n$ if we choose g such that

$$\mathbb{E}[W(y)] < 0.$$

In order to satisfy this inequality, we can search for g among the functions of the form

$$g(y) = \frac{b}{a + \|y\|^2}$$

with an appropriately chosen constants $a \geq 0, b > 0$. Therefore, $W(y)$ can be written as

$$\begin{aligned} W(y) &= -2p\varepsilon^2 \frac{b}{a + \|y\|^2} + 2\varepsilon^2 \sum_{i=1}^p \frac{2by_i^2}{(a + \|y\|^2)^2} + \frac{b^2\|y\|^2}{(a + \|y\|^2)^2} \\ &= \frac{1}{a + \|y\|^2} \left(-2pb\varepsilon^2 + \frac{4b\varepsilon^2\|y\|^2}{a + \|y\|^2} + \frac{b^2\|y\|^2}{(a + \|y\|^2)^2} \right) \\ &\leq \frac{1}{a + \|y\|^2} (-2pb\varepsilon^2 + 4b\varepsilon^2 + b^2) \quad \|y\|^2 \leq a + \|y\|^2 \text{ for } a \geq 0 \\ &= \frac{Q(b)}{a + \|y\|^2}, \quad Q(b) = b^2 - 2pb\varepsilon^2 + 4b\varepsilon^2. \end{aligned}$$

The minimizer in b of quadratic function $Q(b)$ is equal to

$$b_{opt} = \varepsilon^2(p - 2),$$

where the minimum of $W(y)$ satisfies

$$W_{min}(y) \leq -\frac{b_{opt}^2}{a + \|y\|^2} = -\frac{\varepsilon^4(p - 2)^2}{a + \|y\|^2} < 0.$$

Note that when $b \in (b_1, b_2)$, i.e. between the two roots of $Q(b)$

$$b_1 = 0, \quad b_2 = 2\varepsilon^2(p-2)$$

we have $W(y) < 0$, which may lead to other estimators having smaller mean square errors than MLE estimator.

When $a = 0$, the function g and the estimator $\tilde{\mu}_n = (1 - g(y))y$ associated to this choice of g are given by

$$g(y) = \frac{\varepsilon^2(p-2)}{\|y\|^2},$$

and

$$\tilde{\mu}_n = \left(1 - \frac{\varepsilon^2(p-2)}{\|y\|^2}\right)y =: \tilde{\mu}_{JS},$$

respectively. $\tilde{\mu}_{JS}$ is called **James-Stein estimator**. If dimension $p \geq 3$ and the norm $\|y\|^2$ is sufficiently large, multiplication of y by $g(y)$ shrinks the value of y to 0. This is called the **Stein shrinkage**. If $b = b_{opt}$, then

$$W_{min}(y) = -\frac{\varepsilon^4(p-2)^2}{\|y\|^2}.$$

LEMMA 2.6. Let $p \geq 3$. Then, for all $\mu \in \mathbb{R}^p$,

$$0 < \mathbb{E}\left(\frac{1}{\|y\|^2}\right) < \infty.$$

The proof of Lemma 2.2 can be found on Tsybakov's book [Tsy09] (page 158~159). For the function W , Lemma 2.2 implies $-\infty < \mathbb{E}[W(y)] < 0$, provided that $p \geq 3$. Therefore, if $p \geq 3$, the risk of the estimator $\tilde{\mu}_n$ satisfies

$$\mathbb{E}\|\tilde{\mu}_n - \mu\|^2 = p\varepsilon^2 - \mathbb{E}\left(\frac{\varepsilon^4(p-2)^2}{\|y\|^2}\right) < \mathbb{E}\|\hat{\mu}_n - \mu\|^2$$

for all $\mu \in \mathbb{R}^p$.

Besides James-Stein estimator, there are other estimators having smaller mean square errors than MLE $\hat{\mu}_n$.

- Stein estimator: $a = 0, b = \varepsilon^2 p$,

$$\tilde{\mu}_S := \left(1 - \frac{\varepsilon^2 p}{\|y\|^2}\right)y$$

- James-Stein estimator: $c \in (0, 2(p-2))$

$$\tilde{\mu}_{JS}^c := \left(1 - \frac{\varepsilon^2 c}{\|y\|^2}\right)y$$

- Positive part James-Stein estimator:

$$\tilde{\mu}_{JS+} := \left(1 - \frac{\varepsilon^2(p-2)}{\|y\|^2}\right)_+ y$$

- Positive part Stein estimator:

$$\tilde{\mu}_{S+} := \left(1 - \frac{\varepsilon^2 p}{\|y\|^2}\right)_+ y$$

where $(x)_+ = \min(0, x)$. Denote the mean square error by $MSE(\tilde{\mu}) = \mathbb{E}\|\tilde{\mu} - \mu\|^2$, then we have

$$MSE(\tilde{\mu}_{JS+}) < MSE(\tilde{\mu}_{JS}) < MSE(\hat{\mu}_n), \quad MSE(\tilde{\mu}_{S+}) < MSE(\tilde{\mu}_S) < MSE(\hat{\mu}_n).$$

See Efron's Book, Chap 1, Table 1.1.

Another dimension of variation is Shrinkage toward *any vector* rather than the origin.

$$\tilde{\mu}_{\mu_0} = \mu_0 + \left(1 - \frac{\varepsilon^2 c}{\|y\|^2}\right)(y - \mu_0), \quad c \in (0, 2(p-2)).$$

In particular, one may choose $\mu_0 = \bar{y}$ where $\bar{y} = \sum_{i=1}^p y_i/p$.

2.7. Discussion. Stein's phenomenon firstly shows that in high dimensional estimation, shrinkage may lead to better performance than MLE, the sample mean. This opens a new era for modern high dimensional statistics. In fact discussions above study independent random variables in p -dimensional space, concentration of measure tells us some priori knowledge about the estimator distribution – samples are concentrating around certain point. Shrinkage toward such point may naturally lead to better performance.

However, after Stein's phenomenon firstly proposed in 1956, for many years researchers have not found the expected revolution in practice. Mostly because Stein's type estimators are too complicated in real applications and very small gain can be achieved in many cases. Researchers struggle to show real application examples where one can benefit greatly from Stein's estimators. For example, Efron-Morris (1974) showed three examples that JS-estimator significantly improves the multivariate estimation. On other other hand, deeper understanding on Shrinkage-type estimators has been pursued from various aspects in statistics.

The situation changes dramatically when LASSO-type estimators by Tibshirani, also called Basis Pursuit by Donoho et al. are studied around 1996. This brings sparsity and L1-regularization into the central theme of high dimensional statistics and leads to a new type of shrinkage estimator, thresholding. For example,

$$\min_{\tilde{\mu}} I = \min_{\tilde{\mu}} \frac{1}{2}\|\tilde{\mu} - \mu\|^2 + \lambda\|\tilde{\mu}\|_1$$

Subgradients of I over $\tilde{\mu}$ leads to

$$0 \in \partial_{\tilde{\mu}_j} I = (\tilde{\mu}_j - \mu_j) + \lambda \text{sign}(\tilde{\mu}_j) \Rightarrow \tilde{\mu}_j = \text{sign}(\mu_j)(|\mu_j| - \lambda)_+$$

where the set-valued map $\text{sign}(x) = 1$ if $x > 0$, $\text{sign}(x) = -1$ if $x < 0$, and $\text{sign}(x) = [-1, 1]$ if $x = 0$, is the subgradient of absolute function $|x|$. Under this new framework shrinkage estimators achieves a new peak with an ubiquitous spread in data analysis with high dimensionality.

3. Random Matrix Theory and Phase Transitions in PCA

In PCA, one often looks at the eigenvalue plot in an decreasing order as percentage or variations. A large gap in the eigenvalue drops may indicate those top eigenvectors reflect major variation directions, where those small eigenvalues indicate directions due to noise which will vanish when $n \rightarrow \infty$. Is this true in all situations? The situation depend on the following parameter

$$(19) \quad \gamma = \lim_{p,n \rightarrow \infty} \frac{p}{n}.$$

The answer is yes in the classical setting where $\gamma = 0$ governed by the Law of Large Numbers. Unfortunately, in high dimensional statistics with $\gamma > 0$, top eigenvectors of sample covariance matrices might not reflect the subspace of signals. In fact, there is a phase transition for signal identifiability by PCA: below a threshold of signal-noise ratio, PCA will fail with high probability and above that threshold of signal-noise ratio, PCA will approximate the signal subspace with high probability. This will be illustrated by the following simplest rank-1 (spike) signal model, in which a leverage of random matrix theory will shed light on the phase transitions where PCA fails to capture the signal subspace depending on the signal-noise ratio.

3.1. Phase Transitions of PCA in Rank-1 Model. Consider the following rank-1 signal-noise model

$$Y = X + \varepsilon,$$

where

- the signal lies in an one-dimensional subspace $X = \alpha u$ with $\alpha \sim \mathcal{N}(0, \sigma_X^2)$;
- the noise $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2 I_p)$ is i.i.d. Gaussian.

Therefore $Y \sim \mathcal{N}(0, \Sigma)$ where the limiting covariance matrix Σ is a rank-one added by a sparse matrix:

$$\Sigma = \sigma_X^2 uu' + \sigma_\varepsilon^2 I_p.$$

For multi-rank generalizations, please see [KN08].

The whole question in the remaining part of this section is to ask, *can we recover signal direction u from principal component analysis on noisy measurements Y ?*

Define the signal-noise ratio

$$SNR = R = \frac{\sigma_X^2}{\sigma_\varepsilon^2}.$$

For simplicity we assume that $\sigma_\varepsilon^2 = 1$ without loss of generality. We aim to show how SNR affect the result of PCA when p is large. A fundamental result by Johnstone in 2006 [Joh06], or see [NBG10], shows that the primary (largest) eigenvalue of sample covariance matrix satisfies

$$(20) \quad \lambda_{\max}(\widehat{\Sigma}_n) \rightarrow \begin{cases} (1 + \sqrt{\gamma})^2 = b, & \sigma_X^2 \leq \sqrt{\gamma} \\ (1 + \sigma_X^2)(1 + \frac{\gamma}{\sigma_X^2}), & \sigma_X^2 > \sqrt{\gamma} \end{cases}$$

which implies that if signal energy is small, top eigenvalue of sample covariance matrix never pops up from random matrix ones; only if the signal energy is beyond the phase transition threshold $\sqrt{\gamma}$, top eigenvalue can be separated from random matrix eigenvalues. However, even in the latter case it is a biased estimation.

Moreover, the primary eigenvector associated with the largest eigenvalue (principal component) converges to

$$(21) \quad |\langle u, v_{\max} \rangle|^2 \rightarrow \begin{cases} 0 & \sigma_X^2 \leq \sqrt{\gamma} \\ \frac{1 - \frac{\gamma}{\sigma_X^4}}{1 + \frac{\gamma}{\sigma_X^2}}, & \sigma_X^2 > \sqrt{\gamma} \end{cases}$$

which means the same phase transition phenomenon: if signal is of low energy, PCA will tell us nothing about the true signal and the estimated top eigenvector is orthogonal to the true direction u ; if the signal is of high energy, PCA will return a

biased estimation which lies in a cone whose angle with the true signal is no more than

$$\arccos \left(\frac{1 - \frac{\gamma}{\sigma_x^4}}{1 + \frac{\gamma}{\sigma_x^2}} \right).$$

Below we are going to show such results.

3.2. Marčenko-Pastur Law of Sample Covariance Matrix. First of all, we show that even for white noise the sample covariance matrix has its eigenvalue distributed governed by Marčenko-Pastur Law. This is the null distribution describing noise. Let $x_i \sim \mathcal{N}(0, I_p)$ ($i = 1, \dots, n$) and $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{p \times n}$. Then the sample covariance matrix is defined as

$$(22) \quad \widehat{\Sigma}_n = \frac{1}{n} X X^T.$$

Such a random matrix $\widehat{\Sigma}_n$ is called a *Wishart* matrix.

- In classical statistics: when p fixed and $n \rightarrow \infty$, the classical Law of Large Numbers tells us $\widehat{\Sigma}_n \rightarrow I_p$.
- In high dimensional statistics when both n and p grow: $\frac{p}{n} \rightarrow \gamma \neq 0$, the distribution of the eigenvalues of $\widehat{\Sigma}_n$ follows a so called Marčenko-Pastur (MP) distribution [BS10] (Chapter 3),

$$(23) \quad \mu^{MP}(t) = \left(1 - \frac{1}{\gamma} \right) \delta(x) I(\gamma > 1) + \begin{cases} 0 & t \notin [a, b], \\ \frac{\sqrt{(b-t)(t-a)}}{2\pi\gamma t} dt & t \in [a, b], \end{cases}$$

where $a = (1 - \sqrt{\gamma})^2$, $b = (1 + \sqrt{\gamma})^2$. In other words if $\gamma \leq 1$, the distribution has a support on $[a, b]$ and if $\gamma > 1$, it has an additional point mass $1 - 1/\gamma$ at the origin.

Figure 2 illustrates the MP-distribution by MATLAB simulations whose codes can be found below.

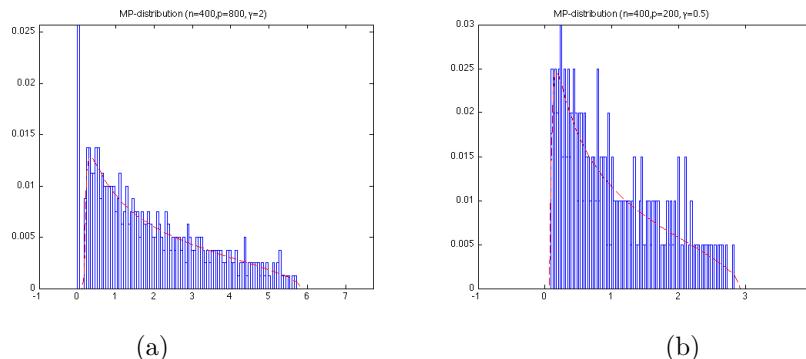


FIGURE 2. (a) Marčenko-Pastur distribution with $\gamma = 2$. (b) Marčenko-Pastur distribution with $\gamma = 0.5$.

```
%Wishart matrix
% S = 1/n*X*X.', X is p-by-n, X_ij i.i.d N(0,1),
% ESD_S converge to M.P. with parameter y = p/n
```

```

y = 2;

a = (1-sqrt(y))^2;
b = (1+sqrt(y))^2;

f_MP = @(t) sqrt(max(b-t, 0).*max(t-a, 0))./(2*pi*y*t); %MP Distribution

%non-zero eigenvalue part
n = 400;
p = n*y;

X = randn(p,n);
S = 1/n*(X*X.');
evals = sort( eig(S), 'descend');

nbin = 100;
[nout, xout] = hist(evals, nbin);
hx = xout(2) - xout(1); % step size, used to compute frequency below
x1 = evals(end) - 1;
x2 = evals(1) + 1; % two end points
xx = x1+hx/2: hx: x2;
fre = f_MP(xx)*hx;

figure,
h = bar(xout, nout/p);
set(h, 'BarWidth', 1, 'FaceColor', 'w', 'EdgeColor', 'b');
hold on;
plot(xx, fre, '--r');

if y > 1 % there are (1-1/y)*p zero eigenvalues
axis([-1 x2+1 0 max(fre)*2]);
end

```

3.3. Characterization of Phase Transitions with RMT. After learning the null distribution of noise is the Marčenko-Pastur Law, now we are ready to come back to the rank-1 spike model.

Following the rank-1 model, consider random vectors $\{Y_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma)$, where $\Sigma = \sigma_x^2 uu^T + \sigma_\varepsilon^2 I_p$ and u is an arbitrarily chosen unit vector ($\|u\|^2 = 1$) showing the signal direction. Define the Signal-Noise-Ratio (SNR) $R = \frac{\sigma_x^2}{\sigma_\varepsilon^2}$. Without loss of generality, we assume $\sigma_\varepsilon^2 = 1$. The covariance matrix Σ thus has a structure that is low-rank plus sparse matrix. The sample covariance matrix is $\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n Y_i Y_i^T = \frac{1}{n} YY^T$ where $Y = [Y_1, \dots, Y_n] \in \mathbb{R}^{p \times n}$. Suppose one of its eigenvalue is λ and the corresponding unit eigenvector is \hat{v} , so $\hat{\Sigma}_n \hat{v} = \lambda \hat{v}$.

First of all, we relate the $\hat{\lambda}$ to the MP distribution by the trick:

$$(24) \quad Z_i = \Sigma^{-\frac{1}{2}} Y_i \rightarrow Z_i \sim N(0, I_p), \text{ where } \Sigma = \sigma_x^2 uu^T + \sigma_\varepsilon^2 I_p = Ruu^T + I_p.$$

Then $S_n = \frac{1}{n} \sum_{i=1}^n Z_i Z_i^T = \frac{1}{n} ZZ^T$ is a Wishart random matrix whose eigenvalues follow the Marčenko-Pastur distribution.

Notice that $\hat{\Sigma}_n = \frac{1}{n} YY^T = \Sigma^{1/2} (\frac{1}{n} ZZ^T) \Sigma^{1/2} = \Sigma^{\frac{1}{2}} S_n \Sigma^{\frac{1}{2}}$ and $(\hat{\lambda}, \hat{v})$ is eigenvalue-eigenvector pair of matrix $\hat{\Sigma}_n$. Therefore

$$(25) \quad \Sigma^{\frac{1}{2}} S_n \Sigma^{\frac{1}{2}} \hat{v} = \hat{\lambda} \hat{v} \Rightarrow S_n \Sigma (\Sigma^{-\frac{1}{2}} \hat{v}) = \hat{\lambda} (\Sigma^{-\frac{1}{2}} \hat{v})$$

In other words, $\hat{\lambda}$ and $\Sigma^{-\frac{1}{2}} \hat{v}$ are the eigenvalue and eigenvector of matrix $S_n \Sigma$. Suppose $c \Sigma^{-\frac{1}{2}} \hat{v} = v$ where the constant c makes v a unit eigenvector and thus satisfies,

$$(26) \quad c^2 = c \hat{v}^T \hat{v} = v^T \Sigma v = v^T (\sigma_x^2 uu^T + \sigma_\varepsilon^2 I_p) v = \sigma_x^2 (u^T v)^2 + \sigma_\varepsilon^2 = R(u^T v)^2 + 1.$$

Now we have,

$$(27) \quad S_n \Sigma v = \hat{\lambda} v.$$

Plugging in the expression of Σ , it gives

$$(28) \quad S_n (\sigma_x^2 uu^T + \sigma_\varepsilon^2 I_p) v = \hat{\lambda} v$$

Rearrange the term with u to one side, we got

$$(29) \quad (\hat{\lambda} I_p - \sigma_\varepsilon^2 S_n) v = \sigma_x^2 S_n u (u^T v)$$

Assuming that $\hat{\lambda} I_p - \sigma_\varepsilon^2 S_n$ is invertable, then multiple its reversion at both sides of the equality, we get,

$$(30) \quad v = \sigma_x^2 \cdot (\hat{\lambda} I_p - \sigma_\varepsilon^2 S_n)^{-1} \cdot S_n u (u^T v).$$

Now we are going to present the estimates on eigenvalue $\hat{\lambda}$ and $|\langle u, v \rangle|$.

3.3.1. Primary Eigenvalue. Multiply (30) by u' at both side,

$$(31) \quad u^T v = \sigma_x^2 \cdot u^T (\hat{\lambda} I_p - \sigma_\varepsilon^2 S_n)^{-1} S_n u \cdot (u^T v)$$

that is, if $u^T v \neq 0$,

$$(32) \quad 1 = \sigma_x^2 \cdot u^T (\hat{\lambda} I_p - \sigma_\varepsilon^2 S_n)^{-1} S_n u$$

Assume that S_n has the eigenvalue decomposition $S_n = W \hat{\Lambda} W^T$, where $\Lambda = \text{diag}(\lambda_i : i = 1, \dots, p)$ ($\lambda_i \geq \lambda_{i+1}$), $WW^T = W^T W = I_p$, with $W = [w_1, w_2, \dots, w_p] \in \mathbb{R}^{p \times p}$ gives an orthonormal basis of eigenvectors. Define $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_n] \in \mathbb{R}^{p \times 1}$, such that $\alpha_i = w_i^T u$, hence $u = \sum_{i=1}^p \alpha_i w_i = W^T \alpha$. Now (32) leads to

$$(33) \quad 1 = \sigma_x^2 \cdot u^T [W (\hat{\lambda} I_p - \sigma_\varepsilon^2 \Lambda)^{-1} W^T] [W \Lambda W^T] u = \sigma_x^2 \cdot \alpha^T (\hat{\lambda} I_p - \sigma_\varepsilon^2 \Lambda)^{-1} \Lambda \alpha$$

which is

$$(34) \quad 1 = \sigma_x^2 \cdot \sum_{i=1}^p \frac{\lambda_i}{\hat{\lambda} - \sigma_\varepsilon^2 \lambda_i} \alpha_i^2$$

where $\sum_{i=1}^p \alpha_i^2 = 1$. Since W consists of a random orthonormal basis on a sphere, α_i will concentrate on its mean $\alpha_i = \frac{1}{\sqrt{q}}$. For large p , $\lambda_i \sim \mu^{MP}(\lambda_i)$ can be thought

sampled from the μ^{MP} and the sum (34) can thus be regarded as the following Monte-Carlo integration with respect to the MP distribution,

$$(35) \quad 1 = \sigma_X^2 \cdot \frac{1}{p} \sum_{i=1}^p \frac{\lambda_i}{\hat{\lambda} - \sigma_\varepsilon^2 \lambda_i} \sim \sigma_X^2 \cdot \int_a^b \frac{t}{\hat{\lambda} - \sigma_\varepsilon^2 t} d\mu^{MP}(t)$$

Since we had assumed without loss of generosity that $\sigma_\varepsilon^2 = 1$, we can compute the integration above using the Stieltjes transform and obtain,

$$(36) \quad 1 = \sigma_X^2 \cdot \int_a^b \frac{t}{\hat{\lambda} - t} \frac{\sqrt{(b-t)(t-a)}}{2\pi\gamma t} dt = \frac{\sigma_X^2}{4\gamma} [2\hat{\lambda} - (a+b) - 2\sqrt{|(\hat{\lambda}-a)(b-\hat{\lambda})|}].$$

For $\hat{\lambda} \geq b$ and $R = \sigma_X^2 \geq \sqrt{\gamma}$, we have

$$\begin{aligned} \therefore \quad 1 &= \frac{\sigma_X^2}{4\gamma} [2\hat{\lambda} - (a+b) - 2\sqrt{(\hat{\lambda}-a)(\hat{\lambda}-b)}], \\ \therefore \quad \hat{\lambda} &= \sigma_X^2 + \frac{\gamma}{\sigma_X^2} + 1 + \gamma = (1 + \sigma_X^2)(1 + \frac{\gamma}{\sigma_X^2}). \end{aligned}$$

More general for $\sigma_\varepsilon^2 \neq 1$, all the equations above is true, except that all the $\hat{\lambda}$ will be replaced by $\frac{\hat{\lambda}}{\sigma_\varepsilon^2}$ and σ_X^2 by signal-noise-ratio $R = \sigma_X^2/\sigma_\varepsilon^2$. Then we get:

$$\hat{\lambda} = (1+R) \left(1 + \frac{\gamma}{R} \right) \sigma_\varepsilon^2.$$

Here we observe the following phase transitions for primary eigenvalue:

- If $\hat{\lambda} \in [a, b]$, then $\hat{\Sigma}_n$ has its primary eigenvalue $\hat{\lambda}$ within $\text{supp}(\mu^{MP})$, so it is undistinguishable from the noise S_n .
- If $\hat{\lambda} \geq b$, PCA will pick up the top eigenvalue as a signal.
- So $\hat{\lambda} = b$ is the phase transition where PCA works to pop up signal rather than noise. Then plugging in $\hat{\lambda} = b$ in (36), we get,

$$(37) \quad 1 = \sigma_X^2 \cdot \frac{1}{4\gamma} [2b - (a+b)] = \frac{\sigma_X^2}{\sqrt{\gamma}} \Leftrightarrow \sigma_X^2 = \sqrt{\frac{p}{n}}$$

Hence, in order to make PCA works, we need to let the signal-noise-ratio $R \geq \sqrt{\frac{p}{n}}$.

3.3.2. Primary Eigenvector. We now study the phase transition of the primary eigenvector. It is convenient to study $|u^T v|^2$ first and then translate back to $|u^T \hat{v}|^2$. From Equation (30), we obtain

$$\begin{aligned} 1 &= v^T v = \sigma_X^4 \cdot v^T u u^T S_n (\lambda I_p - \sigma_\varepsilon^2 S_n)^{-2} S_n u u^T v \\ &= \sigma_X^4 \cdot (|v^T u|) [u^T S_n (\lambda I_p - \sigma_\varepsilon^2 S_n)^{-2} S_n u] (|u^T v|) \end{aligned}$$

which implies that

$$(38) \quad |u^T v|^{-2} = \sigma_X^4 [u^T S_n (\lambda I_p - \sigma_\varepsilon^2 S_n)^{-2} S_n u].$$

Using the same trick as the equation (32), we reach the following Monte-Carlo integration

$$(39) \quad |u^T v|^{-2} = \sigma_X^4 [u^T S_n (\lambda I_p - \sigma_\varepsilon^2 S_n)^{-2} S_n u] \sim \sigma_X^4 \int_a^b \frac{t^2}{(\lambda - \sigma_\varepsilon^2 t)^2} d\mu^{MP}(t)$$

and assume that $\lambda \geq b$, from Stieltjes transform introduced later one can compute the integral as

$$\begin{aligned} |u^T v|^{-2} &= \sigma_X^4 \cdot \int_a^b \frac{t^2}{(\lambda - \sigma_\varepsilon^2 t)^2} d\mu^{MP}(t) \\ &= \frac{\sigma_X^4}{4\gamma} (-4\lambda + (a+b) + 2\sqrt{(\lambda-a)(\lambda-b)} + \frac{\lambda(2\lambda-(a+b))}{\sqrt{(\lambda-a)(\lambda-b)}}) \end{aligned}$$

from which it can be computed that (using $\hat{\lambda} = (1+R)(1+\frac{\gamma}{R})$ obtained above, where $R = \frac{\sigma_X^2}{\sigma_\varepsilon^2}$)

$$|u^T v|^2 = \frac{1 - \frac{\gamma}{R^2}}{1 + \gamma + \frac{2\gamma}{R}}.$$

Now we can compute the inner product of u and \hat{v} that we are really interested in:

$$\begin{aligned} |u^T \hat{v}|^2 &= \left(\frac{1}{c} u^T \Sigma^{\frac{1}{2}} v\right)^2 = \frac{1}{c^2} ((\Sigma^{\frac{1}{2}} u)^T v)^2 \\ &= \frac{1}{c^2} (((Ruu^T + I_p)^{\frac{1}{2}} u)^T v)^2 \\ &\stackrel{*}{=} \frac{1}{c^2} ((\sqrt{(1+R)} u)^T v)^2 \\ &\stackrel{**}{=} \frac{(1+R)(u^T v)^2}{R(u^T v)^2 + 1} \\ &= \frac{1+R - \frac{\gamma}{R} - \frac{\gamma}{R^2}}{1+R + \gamma + \frac{\gamma}{R}} = \frac{1 - \frac{\gamma}{R^2}}{1 + \frac{\gamma}{R}} \end{aligned}$$

where the equality $(*)$ uses $\Sigma^{1/2}u = \sqrt{1+R}u$, and the equality $(**)$ is due to the formula for c^2 (Equation (26) above). Note that this identity holds under the condition that $R \geq \sqrt{\gamma}$ to make the numerator above non-negative.

Therefore if PCA works well and noise doesn't dominate the effect, the inner product $|u^T \hat{v}|$ should be close to 1. Particularly when $\gamma = 0$ we have $|u^T \hat{v}| = 1$ as disclosed by the classical Law of Large Numbers in statistics. On the other hand, from RMT we know that if the top eigenvalue $\hat{\lambda} \in [a, b]$ is overwhelmed in the domain of M. P. distribution, then the primary eigenvector computed from PCA is purely random and $|u^T \hat{v}| = 0$, which means that from \hat{v} we can know nothing about the signal u .

3.4. Stieltjes Transform. Now we present the Stieltjes Transformation of MP-density which has been crucial in computing the integrals above. Define the Stieltjes Transformation of MP-density μ^{MP} to be

$$(40) \quad s(z) := \int_R \frac{1}{t-z} d\mu^{MP}(t), \quad z \in C$$

If $z \in \mathbb{R}$, the transformation is called Hilbert Transformation. Further details can be found in Terry Tao's textbook, Topics on Random Matrix Theory [Tao11], Sec. 2.4.3 (the end of page 169) for the definition of Stieltjes transform of a density $p(t)dt$ on \mathbb{R} .

In [BS10], Lemma 3.11 on page 52 gives the following characterization of $s(z)$ (note that the book contains a typo that $4y\sigma^4$ in numerator should be replaced by

$4yz\sigma^2)$:

$$(41) \quad s(z) = \frac{(1-\gamma)-z+\sqrt{(z-1-\gamma)^2-4\gamma z}}{2\gamma z},$$

which is the largest root of the quadratic equation,

$$(42) \quad \gamma z s(z)^2 + (z - (1 - \gamma))s(z) + 1 = 0 \iff z + \frac{1}{s(z)} = \frac{1}{1 + \gamma s(z)}.$$

From the equation (41), one can take derivative of z on both side to obtain $s'(z)$ in terms of s and z . Using $s(z)$ one can compute the following basic integrals.

LEMMA 3.1. (1)

$$\int_a^b \frac{t}{\lambda-t} \mu^{MP}(t) dt = -\lambda s(\lambda) - 1;$$

(2)

$$\int_a^b \frac{t^2}{(\lambda-t)^2} \mu^{MP}(t) dt = \lambda^2 s'(\lambda) + 2\lambda s(\lambda) + 1$$

PROOF. For convenience, define

$$(43) \quad T(\lambda) := \int_a^b \frac{t}{\lambda-t} \mu^{MP}(t) dt.$$

Note that

$$(44) \quad 1 + T(\lambda) = 1 + \int_a^b \frac{t}{\lambda-t} \mu^{MP}(t) dt = \int_a^b \frac{\lambda-t+t}{\lambda-t} \mu^{MP}(t) dt = -\lambda s(\lambda)$$

which give the first result.

From the definition of $T(\lambda)$, we have

$$(45) \quad \int_a^b \frac{t^2}{(\lambda-t)^2} \mu^{MP}(t) dt = -T(\lambda) - \lambda T'(\lambda).$$

Combined with the first result, we reach the second one. \square

3.5. Further Comments. Random Matrix Theory can only deal with homogeneous Gaussian noise $\sigma_\varepsilon^2 I_p$ here. In practice, Horn's parallel analysis [BE92] is widely used to find the number of principal components or factors using simulations on given data, which has some implementations in R:

<https://cran.r-project.org/web/packages/paran/>

and Matlab

<https://www.mathworks.com/matlabcentral/fileexchange/44996-parallel-analysis--pa--to-for-det>

Moreover, it is still an open problem how to deal with heteroscedastic noise, where Art Owen and Jingshu Wang has some preliminary studies [OW16].

When $\frac{\log(p)}{n} \rightarrow 0$, we need to add more restrictions on $\hat{\Sigma}_n$ in order to estimate it faithfully. There are typically three kinds of restrictions.

- Σ sparse
- Σ^{-1} sparse, also called Precision Matrix
- banded structures (e.g. Toeplitz) on Σ or Σ^{-1}

Recent developments can be found by Bickel, Tony Cai, Tsybakov, Wainwright et al.

For spectral study on random kernel matrices, see El Karoui, Tiefeng Jiang, Xiuyuan Cheng, and Amit Singer et al.

CHAPTER 3

Random Projections and Almost Isometry

1. Introduction

For this class, we introduce Random Projection method which may reduce the dimensionality of n points in \mathbb{R}^p to $k = O(c(\epsilon) \log n)$ at the cost of a uniform metric distortion of at most $\epsilon > 0$, with high probability. The theoretical basis of this method was given as a lemma by Johnson and Lindenstrauss [JL84] in the study of a Lipschitz extension problem. The result has a widespread application in mathematics and computer science. The main application of Johnson-Lindenstrauss Lemma in computer science is high dimensional data compression via random projections [Ach03]. In 2001, Sanjoy Dasgupta and Anupam Gupta [DG03a], gave a simple proof of this theorem using elementary probabilistic techniques in a four-page paper. Below we are going to present a brief proof of Johnson-Lindenstrauss Lemma based on the work of Sanjoy Dasgupta, Anupam Gupta [DG03a], and Dimitris Achlioptas [Ach03].

Recall the problem of MDS: given a set of points $x_i \in \mathbb{R}^p$ ($i = 1, 2, \dots, n$); form a data Matrix $X^{p \times n} = [X_1, X_2 \dots X_n]^T$, when p is large, especially in some cases larger than n , we want to find k -dimensional projection with which pairwise distances of the data point are preserved as well as possible. That is to say, if we know the original pairwise distance $d_{ij} = \|X_i - X_j\|$ or data distances with some disturbance $\tilde{d}_{ij} = \|X_i - X_j\| + \epsilon_{ij}$, we want to find $Y_i \in \mathbb{R}^k$ s.t.:

$$(46) \quad \min \sum_{i,j} (\|Y_i - Y_j\|^2 - d_{ij}^2)^2$$

assuming $\sum_i Y_i = 0$, i.e. putting the origin as data center.

When D is given exactly by the squared distances of points in Euclidean space, classical MDS defines a kernel matrix $B = -\frac{1}{2}H D H$ where $D = (d_{ij}^2)$, $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$, then, the minimization (46) is equivalent to find $Y_i \in \mathbb{R}^k$:

$$(47) \quad \min_{Y \in \mathbb{R}^{k \times n}} \|Y^T Y - B\|_F^2$$

then the row vectors of matrix Y are the eigenvectors (singular vectors) corresponding to k largest eigenvalues (singular values) of B .

The main features of MDS are the following.

- MDS looks for Euclidean embedding of data whose *total* or *average* metric distortion are minimized.
- MDS embedding basis is *adaptive* to the data, namely as a function of data via eigen-decomposition.

Note that distortion measure here amounts to a certain distance between the set of projected points and the original set of points B . Under the Frobenius norm the distortion equals the sum of the squared lengths of these vectors. It is clear that

such vectors captures a significant global property, but it does not offer any local guarantees. Chances are that some points deviate greatly from the original if we only consider the total metric distortion minimization.

What if we want a *uniform* control on metric distortion at every data pair, say

$$(1 - \epsilon) \leq \frac{\|Y_i - Y_j\|^2}{d_{ij}^2} \leq (1 + \epsilon)?$$

Such an embedding is an almost isometry or a Lipschitz mapping from metric space \mathcal{X} to Euclidean space \mathcal{Y} . If \mathcal{X} is an Euclidean space (or more generally Hilbert space), Johnson-Lindenstrauss Lemma tells us that one can take \mathcal{Y} as a subspace of \mathcal{X} of dimension $k = O(c(\epsilon) \log n)$ via random projections to obtain an almost isometry with high probability. As a contrast to MDS, the main features of this approach are the following.

- Almost isometry is achieved with a *uniform* metric distortion bound (*Bi-Lipschitz* bound), with high probability, rather than average metric distortion control;
- The mapping is *universal*, rather than being adaptive to the data.

2. The Johnson-Lindenstrauss Lemma

THEOREM 2.1 (Johnson-Lindenstrauss Lemma). For any $0 < \epsilon < 1$ and any integer n , let k be a positive integer such that

$$k \geq (4 + 2\alpha)(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln n, \quad \alpha > 0.$$

Then for any set V of n points in \mathbb{R}^p , there is a map $f : \mathbb{R}^p \rightarrow \mathbb{R}^k$ such that for all $u, v \in V$

$$(48) \quad (1 - \epsilon) \|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon) \|u - v\|^2$$

Such a f in fact can be found in randomized polynomial time. In fact, inequalities (48) holds with probability at least $1 - 1/n^\alpha$.

REMARK. We have following facts.

- (1) The embedding dimension $k = O(c(\epsilon) \log n)$ which is independent to ambient dimension d and logarithmic to the number of samples n . The independence to d in fact suggests that the Lemma can be generalized to the Hilbert spaces of infinite dimension.
- (2) How to construct the map f ? In fact we can use random projections:

$$Y^{n \times k} = X^{n \times d} R^{d \times k}$$

where the following random matrices R can cater our needs.

- $R = [r_1, \dots, r_k]$ $r_i \in S^{d-1}$ $r_i = (a_1^i, \dots, a_d^i) / \|a^i\|$ $a_k^i \sim N(0, 1)$
- $R = A/\sqrt{k}$ $A_{ij} \sim N(0, 1)$
- $R = A/\sqrt{k}$ $A_{ij} = \begin{cases} 1 & p = 1/2 \\ -1 & p = 1/2 \end{cases}$
- $R = A/\sqrt{k/3}$ $A_{ij} = \begin{cases} 1 & p = 1/6 \\ 0 & p = 2/3 \\ -1 & p = 1/6 \end{cases}$

The proof below actually takes the first form of R as an illustration.

Now we are going to prove Johnson-Lindenstrauss Lemma using a random projection to k -subspace in \mathbb{R}^d . Notice that the distributions of the following two events are identical:

$$\begin{aligned} & \text{unit vector was randomly projected to } k\text{-subspace} \\ \iff & \text{random vector on } S^{d-1} \text{ fixed top-}k \text{ coordinates.} \end{aligned}$$

Based on this observation, we change our target from random k -dimensional projection to random vector on sphere S^{d-1} .

Let $x_i \sim N(0, 1)$ ($i = 1, \dots, p$), and $X = (x_1, \dots, x_p)$, then $Y = X/\|x\| \in S^{p-1}$ is uniformly distributed. Fixing top- k coordinates, we get $z = (x_1, \dots, x_k, 0, \dots, 0)^T/\|x\| \in \mathbb{R}^p$. Let $L = \|z\|^2$ and $\mu := k/p$. Note that $\mathbb{E}\|(x_1, \dots, x_k, 0, \dots, 0)\|^2 = k = \mu \cdot \mathbb{E}\|x\|^2$. The following lemma shows that L is concentrated around μ .

The following lemma is crucial to reach the main theorem.

LEMMA 2.2. let any $k < p$ then we have

(a) if $\beta < 1$ then

$$\text{Prob}[L \leq \beta\mu] \leq \beta^{k/2} \left(1 + \frac{(1-\beta)k}{p-k}\right)^{(p-k)/2} \leq \exp\left(\frac{k}{2}(1-\beta + \ln\beta)\right)$$

(b) if $\beta > 1$ then

$$\text{Prob}[L \geq \beta\mu] \leq \beta^{k/2} \left(1 + \frac{(1-\beta)k}{p-k}\right)^{(p-k)/2} \leq \exp\left(\frac{k}{2}(1-\beta + \ln\beta)\right)$$

Here $\mu = k/p$.

We first show how to use this lemma to prove the main theorem – Johnson-Lindenstrauss lemma.

PROOF OF JOHNSON-LINDENSTRAUSS LEMMA. If $p \leq k$, the theorem is trivial. Otherwise take a random k -dimensional subspace S , and let v'_i be the projection of point $v_i \in V$ into S , then setting $L = \|v'_i - v'_j\|^2$ and $\mu = (k/p)\|v_i - v_j\|^2$ and applying Lemma 2(a), we get that

$$\begin{aligned} \text{Prob}[L \leq (1-\epsilon)\mu] &\leq \exp\left(\frac{k}{2}(1-(1-\epsilon)+\ln(1-\epsilon))\right) \\ &\leq \exp\left(\frac{k}{2}(\epsilon - (\epsilon + \frac{\epsilon^2}{2}))\right), \\ &\quad \text{by } \ln(1-x) \leq -x - x^2/2 \text{ for } 0 \leq x < 1 \\ &= \exp\left(-\frac{k\epsilon^2}{4}\right) \\ &\leq \exp(-(2+\alpha)\ln n), \quad \text{for } k \geq 4(1+\alpha/2)(\epsilon^2/2)^{-1}\ln n \\ &= \frac{1}{n^{2+\alpha}} \end{aligned}$$

$$\begin{aligned}
\text{Prob}[L \geq (1 + \epsilon)\mu] &\leq \exp\left(\frac{k}{2}(1 - (1 + \epsilon) + \ln(1 + \epsilon))\right) \\
&\leq \exp\left(\frac{k}{2}(-\epsilon + (\epsilon - \frac{\epsilon^2}{2} + \frac{\epsilon^3}{3}))\right), \\
&\quad \text{by } \ln(1 + x) \leq x - x^2/2 + x^3/3 \text{ for } x \geq 0 \\
&= \exp\left(-\frac{k}{2}(\epsilon^2/2 - \epsilon^3/3)\right), \\
&\leq \exp(-(2 + \alpha)\ln n), \quad \text{for } k \geq 4(1 + \alpha/2)(\epsilon^2/2 - \epsilon^3/3)^{-1}\ln n \\
&= \frac{1}{n^{2+\alpha}}
\end{aligned}$$

Now set the map $f(x) = \sqrt{\frac{d}{k}}x' = \sqrt{\frac{d}{k}}(x_1, \dots, x_k, 0, \dots, 0)$. By the above calculations, for some fixed pair i, j , the probability that the distortion

$$\frac{\|f(v_i) - f(v_j)\|^2}{\|v_i - v_j\|^2}$$

does not lie in the range $[(1 - \epsilon), (1 + \epsilon)]$ is at most $\frac{2}{n^{(2+\alpha)}}$. Using the trivial union bound with $\binom{n}{2}$ pairs, the chance that some pair of points suffers a large distortion is at most:

$$\binom{n}{2} \frac{2}{n^{(2+\alpha)}} = \frac{1}{n^\alpha} \left(1 - \frac{1}{n}\right) \leq \frac{1}{n^\alpha}.$$

Hence f has the desired properties with probability at least $1 - \frac{1}{n^\alpha}$. This gives us a randomized polynomial time algorithm. \square

Now, it remains to Lemma 2.2.

PROOF OF LEMMA 2.2.

$$\begin{aligned}
\text{Prob}(L \leq \beta\mu) &= \text{Prob}\left(\sum_{i=1}^k x_i^2 \leq \beta\mu(\sum_{i=1}^p x_i^2)\right) \\
&= \text{Prob}\left(\beta\mu \sum_{i=1}^p x_i^2 - \sum_{i=1}^k x_i^2 \geq 0\right) \\
&= \text{Prob}\left[\exp\left(t\beta\mu \sum_{i=1}^p x_i^2 - t \sum_{i=1}^k x_i^2\right) \geq 1\right] \quad (t > 0) \\
&\leq \mathbb{E}\left[\exp\left(t\beta\mu \sum_{i=1}^p x_i^2 - t \sum_{i=1}^k x_i^2\right)\right] \quad (\text{by Markov's inequality}) \\
&= \prod_{i=1}^k \mathbb{E} \exp(t(\beta\mu - 1)x_i^2) \prod_{i=k+1}^p \mathbb{E} \exp(t\beta\mu x_i^2) \\
&= (\mathbb{E} \exp(t(\beta\mu - 1)x^2))^k (\mathbb{E} \exp(t\beta\mu x^2))^{p-k} \\
&= (1 - 2t(\beta\mu - 1))^{-k/2} (1 - 2t\beta\mu)^{-(p-k)/2}
\end{aligned}$$

We use the fact that if $X \sim N(0, 1)$, then $\mathbb{E}[e^{sX^2}] = \frac{1}{\sqrt{(1-2s)}}$, for $-\infty < s < 1/2$.

Now we will refer to last expression as $g(t)$. The last line of derivation gives us the additional constraints that $t\beta\mu \leq 1/2$ and $t(\beta\mu - 1) \leq 1/2$, and so we have $0 < t < 1/(2\beta\mu)$. Now to minimize $g(t)$, which is equivalent to maximize

$$h(t) = 1/g(t) = (1 - 2t(\beta\mu - 1))^{k/2}(1 - 2t\beta\mu)^{(p-k)/2}$$

in the interval $0 < t < 1/(2\beta\mu)$. Setting the derivative $h'(t) = 0$, we get the maximum is achieved at

$$t_0 = \frac{1 - \beta}{2\beta(p - \beta k)}$$

Hence we have

$$h(t_0) = \left(\frac{p - k}{p - k\beta} \right)^{(p-k)/2} \left(\frac{1}{\beta} \right)^{k/2},$$

and this is exactly what we need.

The proof of Lemma 2.2 (b) is almost exactly the same as that of Lemma 2.2 (a). \square

2.1. Conclusion. As we can see, this proof of Lemma is both simple (using just some elementary probabilistic techniques) and elegant. And you may find in the field of machine learning, stochastic method always turns out to be really powerful. The random projection method we approaching today can be used in many fields especially huge dimensions of data is concerned. For one example, in the term document, you may find it really useful for compared with the number of words in the dictionary, the words included in a document is typically sparse (with a few thousands of words) while the dictionary is hugh. Random projections often provide us a useful tool to compress such data without losing much pairwise distance information.

3. Example: MDS in Human Genome Diversity Project

Now consider a SNPs (Single Nucleid Polymorphisms) dataset in Human Genome Diversity Project (HGDP, http://www.cephb.fr/en/hgdp_panel.php) which consists of a data matrix of n -by- p for $n = 1064$ individuals around the world and $p = 644258$ SNPs. Each entry in the matrix has 0, 1, 2, and 9, representing “AA”, “AC”, “CC”, and “missing value”, respectively. After removing 21 rows with all missing values, we are left with a matrix X of size 1043×644258 .

Consider the projection of 1043 persons on the MDS (PCA) coordinates. Let $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ be the centering matrix. Then define

$$K = HXX^TH = U\Lambda U^T$$

which is a positive semi-definite matrix as centered Gram matrix whose eigenvalue decomposition is given by $U\Lambda U^T$. Taking the first two eigenvectors $\sqrt{\lambda_i}u_i$ ($i = 1, \dots, 2$) as the projections of n individuals, Figure 1 gives the projection plot. It is interesting to note that the point cloud data exhibits a continuous trend of human migration in history: origins from Africa, then migrates to the Middle East, followed by one branch to Europe and another branch to Asia, finally spreading into America and Oceania.

One computational concern is that the high dimensionality caused by $p = 644,258$, which is much larger than the number of samples $n = 1043$. However random projections introduced above will provide us an efficient way to compute MDS (PCA) principal components with an almost isometry.

We randomly select (without replacement) $\{n_i, i = 1, \dots, k\}$ from $1, \dots, p$ with equal probability. Let $R \in \mathbb{R}^{k \times p}$ is a Bernoulli random matrix satisfying:

$$R_{ij} = \begin{cases} 1/k & j = n_i, \\ 0 & \text{otherwise.} \end{cases}$$

Now define

$$\tilde{K} = H(XR^T)(RX^T)H$$

whose eigenvectors leads to new principal components of MDS. In the middle and right, Figure 1 plots the such approximate MDS principal components with $k = 5,000$, and $k = 100,000$, respectively. These plots are qualitatively equivalent to the original one.

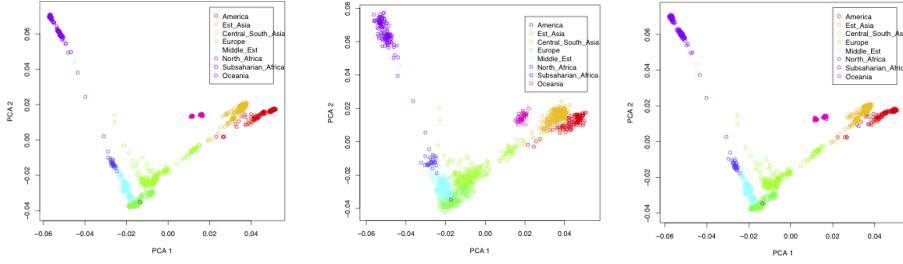


FIGURE 1. (Left) Projection of 1043 individuals on the top 2 MDS principal components. (Middle) MDS computed from 5,000 random projections. (Right) MDS computed from 100,000 random projections. Pictures are due to Qing Wang.

4. Random Projections and Compressed Sensing

There are wide applications of random projections in high dimensional data processing, e.g. [Vem04]. Here we particularly choose a special one, the compressed (or compressive) sensing (CS) where we will use the Johnson-Lindenstrauss Lemma to prove the Restricted Isometry Property (RIP), a crucial result in CS. A reference can be found at [BDDW08].

Compressive sensing can be traced back to 1950s in signal processing in geography. Its modern version appeared in LASSO [Tib96] and BPDN [CDS98], and achieved a highly noticeable status by [CT05, CRT06, CT06].

The basic problem of compressive sensing can be expressed by the following under-determined linear algebra problem. Assume that a signal $x^* \in \mathbb{R}^p$ is sparse with respect to some basis (measurement matrix) $A \in \mathbb{R}^{n \times p}$ or $A \in \mathbb{R}^{n \times p}$ where $n < p$, given measurement $b = Ax^* = Ax^* \in \mathbb{R}^{n^1}$, how can one recover x^* by solving the linear equation system

$$(49) \quad Ax = b?$$

As $n < p$, it is an under-determined problem, whence without further constraint, the problem does not have an unique solution. To overcome this issue, one popular

¹Below we abuse both terms A and Φ while leaving it unique for the future. The noisy version is $b = Ax + \epsilon$ which will be discussed later.

assumption is that the signal x^* is sparse, namely the number of nonzero components $\|x^*\|_0 := \#\{x_i^* \neq 0 : 1 \leq i \leq p\}$ is small compared to the total dimensionality p . Figure 2 gives an illustration of such sparse linear equation problem.

$$\begin{bmatrix} \text{---} \\ n \times I \end{bmatrix} = \begin{bmatrix} \text{---} \\ n \times p \end{bmatrix} \quad \begin{bmatrix} \text{---} \\ p \times 1 \end{bmatrix}$$

FIGURE 2. Illustration of Compressive Sensing (CS). A is a rectangular matrix with more columns than rows. The dark elements represent nonzero elements while the light ones are zeroes. The signal vector x^* , although high dimensional, is sparse.

4.1. Some Sparse Recovery Algorithms. Now we formally give some algorithms to solve the problem. Without loss of generality, we assume each column of design matrix $A = [A_1, \dots, A_p]$ has been standardized, that is, $\|A_j\|_2 = 1$, $j = 1, \dots, p$.

With such a sparse assumption above, a simple idea is to find the sparsest solution satisfying the measurement equation:

$$(P_0) \quad \begin{aligned} & \min && \|x\|_0 \\ & \text{s.t.} && Ax = b. \end{aligned}$$

This is an NP-hard combinatorial optimization problem.

4.1.1. *Basis Pursuit*. A convex relaxation of (50) is called *Basis Pursuit* [CDS98],

$$(P_1) \quad \begin{aligned} & \min && \|x\|_1 := \sum |x_i| \\ & \text{s.t.} && Ax = b. \end{aligned}$$

This is a linear programming problem. Figure 3 shows different projections of a sparse vector x^* under $\|\cdot\|_0$, $\|\cdot\|_1$ and $\|\cdot\|_2$, from which one can see in some cases the convex relaxation (51) does recover the sparse signal solution in (50). Now a natural problem arises, under what conditions the linear programming problem (P_1) has the solution exactly solves (P_0) , i.e. exactly recovers the sparse signal x^* ?

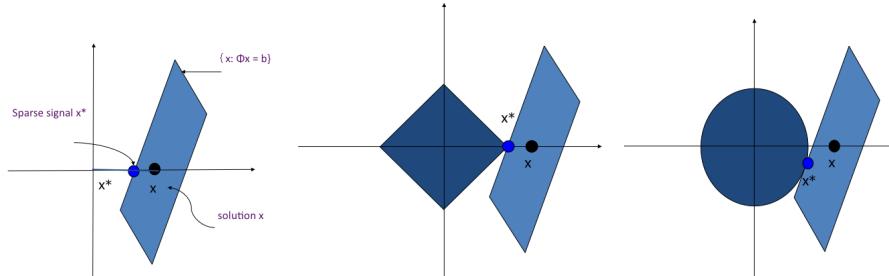


FIGURE 3. Comparison between different projections. Left: projection of x^* under $\|\cdot\|_0$; middle: projection under $\|\cdot\|_1$ which favors sparse solution; right: projection under Euclidean distance.

When measurement noise exists, e.g. $b = Ax + e$ with bound $\|e\|_2$, the following Basis Pursuit De-Noising (BPDN) [CDS98] are used instead

$$(52) \quad \begin{aligned} (BPDN) \quad & \min \|x\|_1 \\ & s.t. \quad \|Ax - b\|_2 \leq \epsilon. \end{aligned}$$

4.1.2. *Orthogonal Matching Pursuit*. Another popular algorithm is proposed by Stephane Mallat and Zhifeng Zhang, 1993 [MZ93].

Algorithm 5: Orthogonal Matching Pursuit (OMP), Mallat-Zhang,1993

Input: A, b
Output: $x \in \mathbb{R}^p$

- 1 initialization: $r_0 = b$, $x_0 = 0$, $S_0 = \emptyset$
- 2 **while** $\|r_t\|_2 \neq 0$ **do**
- 3 $j_t = \arg \max_{1 \leq j \leq p} |\langle A_j, r_{t-1} \rangle|$
- 4 $S_t = S_{t-1} \cup j_t$
- 5 $x_t = \arg \min_{x \in \mathbb{R}^p} \|b - A_{S_t} x\|$
- 6 $r_t = b - Ax_t$
- 7 **end**

Remarks:

1. OMP choose the column of maximal correlation with residue, in fact, it's also the one having the steepest decline in residue, which implies OMP is greedy.
2. In noisy case, we stop algorithm until $\|r_t\| \leq \varepsilon$ for a given ε .
3. It's natural to ask how well OMP can recover x^* , the answer is yes under some conditions we will talk below.

4.1.3. *LASSO*. Least Absolute Shrinkage and Selection Operator [Tib96] solves the following problem for noisy measurement $b = Ax + e$:

$$(53) \quad (LASSO) \quad \min_{x \in \mathbb{R}^p} \|Ax - b\|_2 + \lambda \|x\|_1$$

4.1.4. *Dantzig Selector*. The Dantzig Selector [?] is proposed to deal with noisy measurement $b = Ax + e$:

$$(54) \quad \begin{aligned} & \min \|x\|_1 \\ & s.t. \quad \|A^T(Ax - b)\|_\infty \leq \lambda \end{aligned}$$

For bounded noise $\|e\|_\infty$, the following formulation is used in network analysis [JYLG12]

$$(55) \quad \begin{aligned} & \min \|x\|_1 \\ & s.t. \quad \|Ax - b\|_\infty \leq \epsilon \end{aligned}$$

4.1.5. *Differential Inclusions and Linearized Bregman Iterations*. **to-be-finished...** where

- $\ell(x)$ measures the (empirical) loss of model on a set of samples, such as the mean square loss $\ell(x) = \frac{1}{2n} \|Ax - b\|^2$; and

- $\Omega(x)$ is a sparsity enforcement penalty function, with typical examples including Lasso ($\|x\|_1$), elastic-net ($\alpha\|x\|_1 + (1 - \alpha)\|x\|^2$), group Lasso, and matrix nuclear norm.

The preceding dynamics above are *differential inclusions* due to the inclusive constraint in (??).

The unique feature of (??) includes its dynamics, which follows a regularization path as a family of models at various levels of sparsity or parsimony. [?, **BGOX06**, ?, **BFOS07**] and [**Bur08**] studied it as the *Inverse Scale Space* method in image restoration, and [?] conducted a careful analysis and implementation. The key observation of these works was that large-scale (image) features were recovered before small-scale ones following the dynamics, from which the method derived its name. The discretizations of (??) are known as the Linearized Bregman Iteration (LBI) ([**OBG⁺05**], or equation (3.7) and (5.19-20) of [**YODG08**]), which is developed independently of the continuous dynamics and has been widely used in image processing and compressed sensing. In particular, in terms of matrix completion, [?] studied a discretized version of (??) with a matrix nuclear norm.

Statistical path consistency for (??) was established in [**ORX⁺16**] with its discretized algorithm, the LBI, under sparse linear regression models. This work showed that (??) is able to remove the well-known statistical bias in Lasso estimators [**FL01**], and produce unbiased estimators under nearly the same model selection consistency conditions as Lasso. Furthermore, [**HSXY18**] showed that under a *strictly weaker condition* than generalized Lasso [**LST13**], statistical path consistency could be achieved by (??) equipped with the variable splitting technique. These studies laid down a theoretical foundation for the statistical consistency of regularization paths generated by the solutions of differential inclusion (??). A free R package is released, Libra (Linearized Bregman Algorithms) [**RXY18**]. Another Matlab package² was also released with the Split LBI algorithm [**HSXY16**]. These works fostered various successful applications, such as high dimensional statistics [**XRY18**], computer vision [**FHX⁺16**, **ZSF⁺18**], medical image analysis [**SHYW17**], multimedia [**XXCY16b**], machine learning [**XXCY16a**, **HSXY16**], and AI [**HY18**].

The Linearized Bregman Iterations study the following damped dynamics with $\kappa > 0$ that converges to the (??) as $\kappa \rightarrow \infty$,

Its Euler forward discretization gives where $z_{k+1} = \rho_{k+1} + \frac{x_{k+1}}{\kappa}$, the initial choice $z_0 = \gamma_0 = 0 \in \mathbb{R}^m$ (or small Gaussian), $\beta_0 = 0 \in \mathbb{R}^p$ (or small Gaussian), parameters $\kappa > 0$, $\alpha > 0$, $\nu > 0$, and the proximal map associated with a convex function Ω is defined by

$$\text{prox}_\Omega(z) = \arg \min_x \frac{1}{2} \|z - x\|^2 + \Omega(x).$$

4.2. Uniformly Sparse Recovery Conditions. We are interested in *under which conditions we can recover x^* by those algorithms, and how we can do it when $k = |S| = |\text{supp}(x^*)| \ll n < p$* ?

Now we turn to consider the conditions under which the algorithms above can recover x^* . Below A^* denotes the conjugate of matrix A , which is A^T if A is real.

²<https://github.com/yuany-pku/split-lbi>

a) Uniqueness condition:

$$A_S^* A_S \geq rI, \text{ for some } r > 0.$$

b) Incoherence (Donoho-Huo, 1999): Donoho-Huo [DH01] shows the following sufficient condition

$$\mu = \max_{i \neq j} |\langle A_i, A_j \rangle| < \frac{1}{2k-1},$$

for sparse recovery by BP. This condition is numerically verifiable, so the simplest condition.

c) Irrepresentable condition (Tropp 2004): Joel Tropp shows that under the following condition [Tro04]

$$M =: \|A_{S^c}^* A_S (A_S^* A_S)^{-1}\|_\infty < 1,$$

both OMP and BP recover x^* . This condition is impossible to verify unless the true support set S is known.

d) Restricted-Isometry-Property(R.I.P.)(Candés-Recht-Tao, 2006[CRT06]):

For all k -sparse $x \in \mathbb{R}^p$, $\exists \delta_k \in (0, 1)$, s.t.

$$(1 - \delta_k) \|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_k) \|x\|_2^2.$$

Remarks:

1. Uniqueness condition is a basic one, without which we can't even know which x^* we're going to recover.
2. Irrepresentable condition describe the relevance between A_S and A_{S^c} should be controlled. However, we may regard rows of $A_{S^c}^* A_S (A_S^* A_S)^{-1}$ to be the regression coefficient of $A_j = A_S \beta + \varepsilon$, for $j \in S^c$.
3. In fact, Irrepresentable condition could not be verified before we already have the knowledge about $S = \text{supp}(x^*)$.
4. The only condition easy to check is the incoherence condition, which is stronger than the Irrepresentable condition.
5. The weakest condition, R.I.P. is not easily to be verified. But **Johnson-Lindestrauss Lemma** says some suitable random matrices will satisfy R.I.P. with high probability.

4.2.1. *Incoherence vs. Irrepresentable condition.* Tropp shows that incoherence condition is stronger than the Irrepresentable condition in the following sense:

LEMMA 4.1. (Tropp, 2004 [Tro04])

$$(57) \quad \mu < \frac{1}{2k-1} \Rightarrow M \leq \frac{k\mu}{1 - (k-1)\mu} < 1.$$

Tony Cai et al. shows that the irrepresentable or the incoherence condition is *tight* in the sense that if it fails, there exists data A , x^* , and b such that sparse recovery is not possible.

PROOF OF LEMMA 4.1. First, we have

$$(58) \quad M = \|A_{S^c}^* A_S (A_S^* A_S)^{-1}\|_\infty \leq \|(A_S^* A_S)^{-1}\|_\infty \|A_{S^c}^* A_S\|_\infty.$$

It's easy to verify that

$$(59) \quad \|A_{S^c}^* A_S\|_\infty \leq k\mu.$$

Then we consider $\|(A_S^* A_S)^{-1}\|_\infty$.

Decompose $A_S^* A_S = I_k + \Delta$, then

$$\begin{aligned}
(60) \quad & \max |\Delta_{i,j}| \leq \mu, \quad \text{diag}(\Delta) = \mathbf{0}; \\
& \Rightarrow \|\Delta\|_\infty \leq \frac{k-1}{2k-1} < 1; \\
& \Rightarrow (A_S^* A_S)^{-1} = (I_k + \Delta)^{-1} = \sum_{j=0}^{\infty} (-\Delta)^j; \\
& \Rightarrow \|(A_S^* A_S)^{-1}\|_\infty = \left\| \sum_{j=0}^{\infty} (-\Delta)^j \right\|_\infty \leq \sum_{j=0}^{\infty} \|\Delta\|_\infty^j = \frac{1}{1 - \|\Delta\|_\infty} \leq \frac{1}{1 - (k-1)\mu}.
\end{aligned}$$

Thus, we reach our conclusion

$$(61) \quad M \leq \frac{k\mu}{1 - (k-1)\mu}.$$

□

[DH01] shows that as long as sparsity of x^* satisfies

$$\|x^*\|_0 = |T| < \frac{1 + \frac{1}{\mu(A)}}{2}$$

which is later improved by [EB01] to be

$$\|x^*\|_0 = |T| < \frac{\sqrt{2} - \frac{1}{2}}{\mu(A)},$$

then P_1 recovers x^* .

THEOREM 4.2. (Tropp, 2004) Under uniqueness and Irrepresentable conditions, OMP and BP recovers x^* .

PROOF OF THEOREM 4.2. (I) OMP recovers x^* .

The key to the proof is to show that at each step $t \leq k$, OMP selects atom from S rather than S^c . Then we only need to examine

$$(62) \quad \rho(r_t) = \frac{\|A_{S^c}^* r_t\|_\infty}{\|A_S^* r_t\|_\infty} < 1.$$

In noise-free case,

$$(63) \quad \left. \begin{aligned} b &= Ax^* \in \text{im}(A_S) \\ r_t &= b - Ax_t \in \text{im}(A_S) \end{aligned} \right\} \Rightarrow r_t \in \text{im}(A_S).$$

$P_S = A_S (A_S^* A_S)^{-1} A_S^*$ is the projection operator onto $\text{im}(A_S)$, thus we have $r_t = P_S r_t$. Hence,

$$(64) \quad \rho(r_t) = \frac{\|A_{S^c}^* (P_S r_t)\|_\infty}{\|A_S^* r_t\|_\infty} = \frac{\|A_{S^c}^* A_S (A_S^* A_S)^{-1} A_S^* r_t\|_\infty}{\|A_S^* r_t\|_\infty} \leq \|A_{S^c}^* A_S (A_S^* A_S)^{-1}\|_\infty < 1.$$

(II) BP recovers x^* .

Assume $\hat{x} \neq x^*$ solves

$$(65) \quad P_1 : \min \|x\|_1, \quad \text{s.t.} \quad Ax = b.$$

Denote $\hat{S} = \text{supp}(\hat{x})$ and $\hat{S} \setminus S \neq \emptyset$. We have

$$\begin{aligned}
(66) \quad \|x^*\|_1 &= \|(A_S^* A_S)^{-1} A_S^* b\|_1 \\
&= \|(A_S^* A_S)^{-1} A_S^* A_{\hat{S}} \hat{x}_{\hat{S}}\|_1 \quad (A\hat{x} = b) \\
&= \|(A_S^* A_S)^{-1} A_S^* A_S \hat{x}_S + (A_S^* A_S)^{-1} A_S^* A_{\hat{S} \setminus S} \hat{x}_{\hat{S} \setminus S}\|_1 \quad (\hat{x}_{\hat{S}} = \hat{x}_S + \hat{x}_{\hat{S} \setminus S}) \\
&< \|\hat{x}_S\|_1 + \|\hat{x}_{\hat{S} \setminus S}\|_1 = \|\hat{x}_{\hat{S}}\|_1,
\end{aligned}$$

which is a contradiction. \square

4.2.2. RIP and Random Projections. [AC09] shows that incoherence conditions implies RIP, whence RIP is a weaker condition. Under RIP condition, uniqueness of P_0 and P_1 can be guaranteed for all k -sparse signals, often called *uniform exact recovery* [Can08].

THEOREM 4.3. The following holds for all k -sparse x^* satisfying $Ax^* = b$.

- (1) If $\delta_{2k} < 1$, then problem P_0 has a unique solution x^* ;
- (2) If $\delta_{2k} < \sqrt{2} - 1$, then the solution of P_1 (51) has a unique solution x^* , i.e. recovers the original sparse signal x^* .

The first condition³ is nothing but every $2k$ -columns of A are linearly dependent. To see the first condition, assume by contradiction that there is another k -sparse solution of P_0 , x' . Then by $Ay = 0$ and $y = x^* - x'$ is $2k$ -sparse. If $y \neq 0$, it violates $\delta_{2k} < 1$ such that $0 = \|Ay\| \geq (1 - \delta_{2k})\|y\| > 0$. Hence one must have $y = 0$, i.e. $x^* = x'$ which proves the uniqueness of P_0 . The proof of the second condition can be found in [Can08].

RIP conditions also lead to upper bounds between solutions above and the true sparse signal x^* . For example, in the case of BPDN the following result holds [Can08].

THEOREM 4.4. Suppose that $\|e\|_2 \leq \epsilon$. If $\delta_{2k} < \sqrt{2} - 1$, then

$$\|\hat{x} - x^*\|_2 \leq C_1 k^{-1/2} \sigma_k^1(x^*) + C_2 \epsilon,$$

where \hat{x} is the solution of BPDN and

$$\sigma_k^1(x^*) = \min_{\text{supp}(y) \leq k} \|x^* - y\|_1$$

is the best k -term approximation error in l_1 of x^* .

How to find matrices satisfying RIP? Equipped with Johnson-Lindenstrauss Lemma, one can construct such matrices by random projections with high probability [BDDW08].

Recall that in the Johnson-Lindenstrauss Lemma, a random matrix $A \in \mathbb{R}^{n \times p}$ with each element is i.i.d. according to some distribution satisfying certain bounded moment conditions, e.g. $A_{ij} \sim \mathcal{N}(0, 1)$. The key step to establish Johnson-Lindenstrauss Lemma is the following fact

$$(67) \quad \Pr(|\|Ax\|_2^2 - \|x\|_2^2| \geq \epsilon \|x\|_2^2) \leq 2e^{-nc_0(\epsilon)}.$$

³The necessity of the first condition fails. As pointed to me by Mr. Kaizheng Wang, a counter example can be constructed as follows: Let $A = [1, 1, 1, 0; 1, 1, 0, 1]$, $x^* = [0, 0, 1, 0]^T$, $b = [1, 0]^T$, $x = [1, -1, 0, 0]^T$, $k = 1$. Then x^* is the unique k -sparse solution to $Ax^* = b$. On the other hand, x is $2k$ -sparse, but $Ax = 0$. Hence dependence of columns in A implies that $\delta_{2k} \geq 1$ which disproves necessity of $\delta_{2k} < 1$.

With this one can establish a bound on the action of A on k -sparse x by an union bound via covering numbers of k -sparse signals.

LEMMA 4.5. Let $A \in \mathbb{R}^{n \times p}$ be a random matrix satisfying the concentration inequality (67). Then for any $\delta \in (0, 1)$ and any set all T with $|T| = k < n$, the following holds

$$(68) \quad (1 - \delta)\|x\|_2 \leq \|Ax\|_2 \leq (1 + \delta)\|x\|_2$$

for all x whose support is contained in T , with probability at least

$$(69) \quad 1 - 2 \left(\frac{12}{\delta} \right)^k e^{-c_0(\delta/2)n}.$$

PROOF. It suffices to prove the results when $\|x\|_2 = 1$ as A is linear. Let $X_T := \{x : \text{supp}(x) = T, \|x\|_2 = 1\}$. We first choose Q_T , a $\delta/4$ -cover of X_T , such that for every $x \in X_T$ there exists $q \in Q_T$ satisfying $\|q - x\|_2 \leq \delta/4$. Since X_T has dimension at most k , it is well-known from covering numbers that the capacity $\#(Q_T) \leq (12/\delta)^k$. Now we are going to apply the union bound of (67) to the set Q_T with $\epsilon = \delta/2$. For each $q \in Q_T$, with probability at most $2e^{-c_0(\delta/2)n}$, $\|Aq\|_2^2 - \|q\|_2^2 \geq \delta/2\|q\|_2^2$. Hence for all $q \in Q_T$, the same bound holds with probability at most

$$2\#(Q_T)e^{-c_0(\delta/2)n} \leq 2 \left(\frac{12}{\delta} \right)^k e^{-c_0(\delta/2)n}.$$

Now we define α to be the smallest constant such that

$$\|Ax\|_2 \leq (1 + \alpha)\|x\|_2, \quad \text{for all } x \in X_T.$$

We can show that $\alpha \leq \delta$ with the same probability. For this, pick up a $q \in Q_T$ such that $\|q - x\|_2 \leq \delta/4$, whence by the triangle inequality

$$\|Ax\|_2 \leq \|Aq\|_2 + \|A(x - q)\|_2 \leq 1 + \delta/2 + (1 + \alpha)\delta/4.$$

This implies that $\alpha \leq \delta/2 + (1 + \alpha)\delta/4$, whence $\alpha \leq 3\delta/4/(1 - \delta/4) \leq \delta$. This gives the upper bound. The lower bound also follows this since

$$\|Ax\|_2 \geq \|Aq\|_2 - \|A(x - q)\|_2 \geq 1 - \delta/2 - (1 + \delta)\delta/4 \geq 1 - \delta,$$

which completes the proof. \square

With this lemma, note that there are at most $\binom{p}{k}$ subspaces of k -sparse, an union bound leads to the following result for RIP.

THEOREM 4.6. Let $A \in \mathbb{R}^{n \times p}$ be a random matrix satisfying the concentration inequality (67) and $\delta \in (0, 1)$. There exists $c_1, c_2 > 0$ such that if

$$k \leq c_1 \frac{n}{\log(p/k)}$$

the following RIP holds for all k -sparse x ,

$$(1 - \delta_k)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_k)\|x\|_2^2$$

with probability at least $1 - 2e^{-c_2 n}$.

PROOF. For each of k -sparse signal (X_T) , RIP fails with probability at most

$$2 \left(\frac{12}{\delta} \right)^k e^{-c_0(\delta/2)n}.$$

There are $\binom{p}{k} \leq (ep/k)^k$ such subspaces. Hence, RIP fails with probability at most

$$2 \left(\frac{ep}{k} \right)^k \left(\frac{12}{\delta} \right)^2 e^{-c_0(\delta/2)n} = 2e^{-c_0(\delta/2)n + k[\log(ep/k) + \log(12/\delta)]}.$$

Thus for a fixed $c_1 > 0$, whenever $k \leq c_1 n / \log(p/k)$, the exponent above will be $\leq -c_2 n$ provided that $c_2 \leq c_0(\delta/2) - c_1(1 + (1 + \log(12/\delta)) / \log(p/k))$. c_2 can be always chosen to be > 0 if $c_1 > 0$ is small enough. This leads to the results. \square

Another use of random projections (random matrices) can be found in Robust Principal Component Analysis (RPCA) in the next chapter.

CHAPTER 4

Generalized PCA/MDS via SDP Relaxations

1. Introduction of SDP with a Comparison to LP

Here we will give a short note on Semidefinite Programming (SDP) formulation of Robust PCA, Sparse PCA, MDS with uncertainty, and Maximal Variance Unfolding, etc. First of all, we give a short introduction to SDP based on a parallel comparison with LP.

Semi-definite programming (SDP) involves linear objective functions and linear (in)equalities constraint with respect to variables as positive semi-definite matrices. SDP is a generalization of linear programming (LP) by replacing nonnegative variables with positive semi-definite matrices. We will give a brief introduction of SDP through a comparison with LP.

LP (Linear Programming): for $x \in \mathbb{R}^n$ and $c \in \mathbb{R}^n$,

$$(70) \quad \begin{aligned} & \min && c^T x \\ & s.t. && Ax = b \\ & && x \geq 0 \end{aligned}$$

This is the primal linear programming problem.

In SDP, the inner product between vectors $c^T x$ in LP will change to Hadamard inner product (denoted by \bullet) between matrices.

SDP (Semi-definite Programming): for $X, C \in \mathbb{R}^{n \times n}$

$$(71) \quad \begin{aligned} & \min && C \bullet X = \sum_{i,j} c_{ij} X_{ij} \\ & s.t. && A_i \bullet X = b_i, \quad \text{for } i = 1, \dots, m \\ & && X \succeq 0 \end{aligned}$$

Linear programming has a dual problem via the Lagrangian. The Lagrangian of the primal problem is

$$\max_{\mu \geq 0, y} \min_x L_{x;y,\mu} = c^T x + y^T (b - Ax) - \mu^T x$$

which implies that

$$\begin{aligned} \frac{\partial L}{\partial x} &= c - A^T y - \mu = 0 \\ \iff c - A^T y &= \mu \geq 0 \\ \implies \max_{\mu \geq 0, y} L &= -y^T b \end{aligned}$$

which leads to the following dual problem.

LD (Dual Linear Programming):

$$(72) \quad \begin{aligned} & \max \quad b^T y \\ & s.t. \quad \mu = c - A^T y \geq 0 \end{aligned}$$

In a similar manner, for SDP's dual form, we have the following.

SDD (Dual Semi-definite Programming):

$$(73) \quad \begin{aligned} & \max \quad b^T y \\ & s.t. \quad S = C - \sum_{i=1}^m A_i y_i \succeq 0 =: C - A^T \otimes y \end{aligned}$$

where

$$A = \begin{bmatrix} A_1 \\ \vdots \\ A_m \end{bmatrix}$$

and

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

1.1. Duality of SDP. Define the feasible set of primal and dual problems are $\mathbb{F}_p = \{X \succeq 0; A_i \bullet X = b_i\}$ and $\mathbb{F}_d = \{(y, S) : S = C - \sum_i y_i A_i \succeq 0\}$, respectively. Similar to linear programming, semi-definite programming also has properties of weak and strong duality. The weak duality says that the primal value is always an upper bound of dual value. The strong duality says that the existence of an interior point ensures the vanishing duality gap between primal value and dual value, as well as the complementary conditions. In this case, to check the optimality of a primal variable, it suffices to find a dual variable which meets the complementary condition with the primal. This is often called the *witness* method.

For more reference on duality of SDP, see e.g. [Ali95].

THEOREM 1.1 (Weak Duality of SDP). If $\mathbb{F}_p \neq \emptyset, \mathbb{F}_d \neq \emptyset$, We have $C \bullet X \geq b^T y$, for $\forall X \in \mathbb{F}_p$ and $\forall (y, S) \in \mathbb{F}_d$.

THEOREM 1.2 (Strong Duality SDP). Assume the following hold,

- (1) $\mathbb{F}_p \neq \emptyset, \mathbb{F}_d \neq \emptyset$;
- (2) At least one feasible set has an interior.

Then X^* is optimal iff

- (1) $X^* \in \mathbb{F}_p$
- (2) $\exists (y^*, S^*) \in \mathbb{F}_d$

s.t. $C \bullet X^* = b^T y^*$ or $X^* S^* = 0$ (note: in matrix product)

In other words, the existence of an interior solution implies the complementary condition of optimal solutions. Under the complementary condition, we have

$$\text{rank}(X^*) + \text{rank}(S^*) \leq n$$

for every optimal primal X^* and dual S^* .

2. Robust PCA

Let $X \in \mathbb{R}^{p \times n}$ be a data matrix. Classical PCA tries to find

$$(74) \quad \begin{aligned} & \min \|X - L\| \\ & \text{s.t. } \text{rank}(L) \leq k \end{aligned}$$

where the norm here is any unitary invariant matrix norms, e.g. Schatten's p -norm $\|M\|_p = (\sum_i \sigma_i(M)^p)^{1/p}$ ($p \geq 1$) when M admits the Singular Value Decomposition (SVD) $M = USV^T$ with $S = \text{diag}(\sigma_1, \dots, \sigma_k, \dots)$ ($p = 2$ is the Frobenius norm, $p = 1$ is the nuclear norm, and $p = \infty$ gives spectral norm). SVD provides a solution with $L = \sum_{i \leq k} \sigma_i u_i v_i^T$ where $X = \sum_i \sigma_i u_i v_i^T$ ($\sigma_1 \geq \sigma_2 \geq \dots$). In other words, classical PCA looks for decomposition

$$X = L + E$$

where the error matrix E has small a Frobenius norm which usually is the case for Gaussian noise. However, if some outliers exists, i.e. there are a small amount of sample points which are largely deviated from the main population of samples, the classical PCA is well-known very sensitive to such outliers.

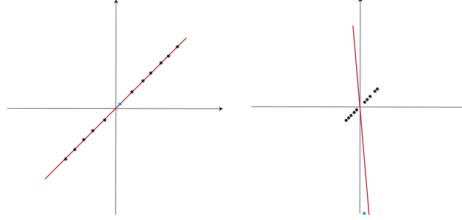


FIGURE 1. Classical PCA is sensitive to outliers

To address this issue, Robust PCA looks for the following decomposition instead

$$X = L + S$$

where

- L is a low rank matrix;
- S is a sparse matrix.

EXAMPLE 4. In the spike signal model, $X = \alpha u + \sigma_\epsilon \epsilon$, where $\alpha \sim \mathcal{N}(0, \sigma_u^2)$ and $\epsilon \sim \mathcal{N}(0, I_p)$. X is thus subject to the following normal distribution $\mathcal{N}(0, \Sigma)$ where $\Sigma = \sigma_u^2 uu^T + \sigma_\epsilon^2 I$. So $\Sigma = L + S$ has such a rank-sparsity structure with $L = \sigma_u^2 uu^T$ and $S = \sigma_\epsilon^2 I$.

EXAMPLE 5. Let $X = [x_1, \dots, x_p]^T \sim \mathcal{N}(0, \Sigma)$ be multivariate Gaussian random variables. The following characterization [CPW12] holds

x_i and x_j are conditionally independent given other variables

$$\Leftrightarrow (\Sigma^{-1})_{ij} = 0$$

We denote it by $x_i \perp x_j | x_k (k \notin \{i, j\})$. Let $G = (V, E)$ be a undirected graph where V represent p random variables and $(i, j) \in E \Leftrightarrow x_i \perp x_j | x_k (k \notin \{i, j\})$. G is called a (Gaussian) graphical model of X .

Divide the random variables into observed and hidden (a few) variables $X = (X_o, X_h)^T$ (in semi-supervised learning, unlabeled and labeled, respectively) and

$$\Sigma = \begin{bmatrix} \Sigma_{oo} & \Sigma_{oh} \\ \Sigma_{ho} & \Sigma_{hh} \end{bmatrix} \quad \text{and} \quad Q = \Sigma^{-1} = \begin{bmatrix} Q_{oo} & Q_{oh} \\ Q_{ho} & Q_{hh} \end{bmatrix}$$

The following Schur Complement equation holds for covariance matrix of observed variables

$$\Sigma_{oo}^{-1} = Q_{oo} + Q_{oh}Q_{hh}^{-1}Q_{ho}.$$

Note that

- Observable variables are often conditional independent given hidden variables, so Q_{oo} is expected to be *sparse*;
- Hidden variables are of small number, so $Q_{oh}Q_{hh}^{-1}Q_{ho}$ is of *low-rank*.

In semi-supervised learning, the labeled points are of small number, and the unlabeled points should be as much conditional independent as possible to each other given labeled points. This implies that the labels should be placed on those most “influential” points.



FIGURE 2. Surveillance video as low rank plus sparse matrices:
Left = low rank (middle) + sparse (right) [CLMW09]

EXAMPLE 6 (Surveillance Video Decomposition). Figure 2 gives an example of low rank vs. sparse decomposition in surveillance video. On the left column, surveillance video of a movie theatre records a great amount of images with the same background and the various walking customers. If we vectorize these images (each image as a vector) to form a matrix, the background image leads to a rank-1 part and the occasional walking customers contribute to the sparse part.

More examples can be found at [**CLMW09**, **CSPW11**, **CPW12**].

In Robust PCA the purpose is to solve

$$(75) \quad \begin{aligned} & \min \quad \|X - L\|_0 \\ & \text{s.t.} \quad \text{rank}(L) \leq k \end{aligned}$$

where $\|A\|_0 = \#\{A_{ij} \neq 0\}$. However both the objective function and the constraint are non-convex, whence it is NP-hard to solve in general.

The simplest convexification leads to a Semi-definite relaxation:

$$\begin{aligned} \|S\|_0 &:= \#\{S_{ij} \neq 0\} \Rightarrow \|S\|_1 \\ \text{rank}(L) &:= \#\{\sigma_i(L) \neq 0\} \Rightarrow \|L\|_* = \sum_i \sigma_i(L), \end{aligned}$$

where $\|L\|_*$ is called the *nuclear norm* of L , which has a semi-definite representation

$$\begin{aligned} \|L\|_* &= \min \quad \frac{1}{2}(\text{trace}(W_1) + \text{trace}(W_2)) \\ \text{s.t.} \quad & \begin{bmatrix} W_1 & L \\ L^T & W_2 \end{bmatrix} \succeq 0. \end{aligned}$$

With these, the relaxed Robust PCA problem can be solved by the following semi-definite programming (SDP).

$$(76) \quad \begin{aligned} & \min \quad \frac{1}{2}(\text{trace}(W_1) + \text{trace}(W_2)) + \lambda \|S\|_1 \\ & \text{s.t.} \quad L_{ij} + S_{ij} = X_{ij}, \quad (i, j) \in E \\ & \quad \begin{bmatrix} W_1 & L \\ L^T & W_2 \end{bmatrix} \succeq 0 \end{aligned}$$

The following Matlab codes realized the SDP algorithm above by CVX (<http://cvxr.com/cvx>).

```
% Construct a random 20-by-20 Gaussian matrix and construct a rank-1
% matrix using its top-1 singular vectors
R = randn(20,20);
[U,S,V] = svds(R,3);
A = U(:,1)*V(:,1)';

% Construct a 90% uniformly sparse matrix
E0 = rand(20);
E = 1*abs(E0>0.9);

X = A + E;

% Choose the regularization parameter
lambda = 0.25;

% Solve the SDP by calling cvx toolbox
if exist('cvx_setup.m','file'),
    cd /matlab_tools/cvx/
    cvx_setup
end
```

```

cvx_begin
    variable L(20,20);
    variable S(20,20);
    variable W1(20,20);
    variable W2(20,20);
    variable Y(40,40) symmetric;
    Y == semidefinite(40);
    minimize(.5*trace(W1)+0.5*trace(W2)+lambda*sum(sum(abs(S))));
    subject to
        L + S >= X-1e-5;
        L + S <= X + 1e-5;
        Y == [W1, L';L W2];
cvx_end

% The difference between sparse solution S and E
disp('S-E\infty')
norm(S-E,'inf')

% The difference between the low rank solution L and A
disp('A-L')
norm(A-L)

```

Typically CVX only solves SDP problem of small sizes (say matrices of size less than 100). Specific matlab tools have been developed to solve large scale RPCA, which can be found at <http://perception.cs1.uiuc.edu/matrix-rank/home.html>.

Some theory based on convex geometry can be found in [**CSPW11**, **CRPW12**]. Besides the SDP approach, some other developments can be found in l_p distance [**LZ11**] and Tyler's M-estimator [**Zha16**, **ZCS14**].

3. Exact Recovery Conditions for RPCA

A fundamental question about Robust PCA is: given $X = L_0 + S_0$ with low-rank L and sparse S , under what conditions that one can recover X by solving SDP in (76)?

It is necessary to assume that

- the low-rank matrix L_0 can not be sparse;
- the sparse matrix S_0 can not be of low-rank.

Such an assumption can be characterized using the following algebraic language.

Define

$$T(L_0) = \{UA^T + BV^T : \forall A, B \in \mathbb{R}^{n \times p}, L_0 = USV^T\}$$

which is the tangent space at L_0 varying in the same column and row spaces of L_0 , and

$$\Omega(S_0) = \{S : \text{supp}(S) \subseteq \text{supp}(S_0)\},$$

which is the tangent space at S_0 varying within the same support of S_0 . The assumptions above are equivalent to say that tangent spaces $T(L_0)$ and $\Omega(S_0)$ are transversal with only intersection at 0,

$$\text{Transversality: } T(L_0) \cap \Omega(S_0) = \{0\}.$$

The following two incoherence constants measure the “diffusive behaviours” of sparse (low-rank) matrices onto low-rank (sparse) opponents.

$$\mu(S_0) = \max_{S \in \Omega(S_0), \|S\|_\infty \leq 1} \|S\|_2$$

$$\xi(L_0) = \max_{L \in T(L_0), \|L\|_2 l \leq 1} \|L\|_\infty$$

[CSPW11] shows the following uncertainty principle, for any matrix M , $\mu(M) \cdot \xi(M) \geq 1$. Therefore a sufficient condition holds,

$$\mu(S_0) \cdot \xi(L_0) < 1, \Rightarrow T(L_0) \cap \Omega(S_0) = \{0\}.$$

Moreover, [CSPW11] shows the following deterministic recovery conditions by SDP

$$\mu(S_0) \cdot \xi(L_0) < 1/6, \Rightarrow \text{SDP recovers } L_0 \text{ and } S_0.$$

Probabilistic recovery conditions are given earlier in [CR09]. First of all we need some incoherence conditions for the identifiability. Assume that $L_0 \in \mathbb{R}^{n \times n} = U\Sigma V^T$ and $r = \text{rank}(L_0)$.

Incoherence condition [CR09]: there exists a $\mu \geq 1$ such that for all $e_i = (0, \dots, 0, 1, 0, \dots, 0)^T$,

$$\|U^T e_i\|^2 \leq \frac{\mu r}{n}, \quad \|V^T e_i\|^2 \leq \frac{\mu r}{n},$$

and

$$|UV^T|_{ij}^2 \leq \frac{\mu r}{n^2}.$$

These conditions, roughly speaking, ensure that the singular vectors are not sparse, i.e. well-spread over all coordinates and won’t concentrate on some coordinates. The incoherence condition holds if $|U_{ij}|^2 \vee |V_{ij}|^2 \leq \mu/n$. In fact, if U represent random projections to r -dimensional subspaces with $r \geq \log n$, we have $\max_i \|U^T e_i\|^2 \asymp r/n$.

To meet the second condition, we simply assume that the sparsity pattern of S_0 is uniformly random.

THEOREM 3.1. Assume the following holds,

- (1) L_0 is n -by- n with $\text{rank}(L_0) \leq \rho_r n \mu^{-1} (\log n)^{-2}$,
- (2) S_0 is uniformly sparse of cardinality $m \leq \rho_s n^2$.

Then with probability $1 - O(n^{-10})$, (76) with $\lambda = 1/\sqrt{n}$ is exact, i.e. its solution $\hat{L} = L_0$ and $\hat{S} = S_0$.

Note that if L_0 is a rectangular matrix of $n_1 \times n_2$, the same holds with $\lambda = 1/\sqrt{(\max n_1, n_2)}$. The result can be generalized to $1 - O(n^{-\beta})$ for $\beta > 0$. Extensions and improvements of these results to incomplete measurements can be found in [CT10, Gro11] etc., which solves the following SDP problem.

$$(77) \quad \begin{aligned} \min \quad & \|L\|_* + \lambda \|S\|_1 \\ \text{s.t.} \quad & L_{ij} + S_{ij} = X_{ij}, \quad (i, j) \in \Omega_{obs}. \end{aligned}$$

THEOREM 3.2. Assume the following holds,

- (1) L_0 is n -by- n with $\text{rank}(L_0) \leq \rho_r n \mu^{-1} (\log n)^{-2}$,
- (2) Ω_{obs} is a uniform random set of size $m = 0.1n^2$,
- (3) each observed entry is corrupted with probability $\tau \leq \tau_s$.

Then with probability $1 - O(n^{-10})$, (76) with $\lambda = 1/\sqrt{0.1n}$ is exact, i.e. its solution $\hat{L} = L_0$. The same conclusion holds for rectangular matrices with $\lambda = 1/\sqrt{\max \dim}$.

All these results hold irrespective to the magnitudes of L_0 and S_0 .

When there are no sparse perturbation in optimization problem (77), the problem becomes the classical Matrix Completion problem with uniformly random sampling:

$$(78) \quad \begin{aligned} & \min \quad \|L\|_* \\ & s.t. \quad L_{ij} = L_{ij}^0, \quad (i, j) \in \Omega_{obs}. \end{aligned}$$

Assumed the same condition as before, [CT10] gives the following result: solution to SDP (78) is exact with probability at least $1 - n^{-10}$ if $m \geq \mu nr \log^a n$ where $a \leq 6$, which can be improved by [Gro11] to be near-optimal

$$m \geq \mu nr \log^2 n.$$

3.1. Phase Transitions. Take $L_0 = UV^T$ as a product of $n \times r$ i.i.d. $\mathcal{N}(0, 1)$ random matrices. Figure 3 shows the phase transitions of successful recovery probability over sparsity ratio $\rho_s = m/n^2$ and low rank ratio r/n . White color indicates the probability equals to 1 and black color corresponds to the probability being 0. A sharp phase transition curve can be seen in the pictures. (a) and (b) respectively use random signs and coherent signs in sparse perturbation, where (c) is purely matrix completion with no perturbation. Increasing successful recovery can be seen from (a) to (c).

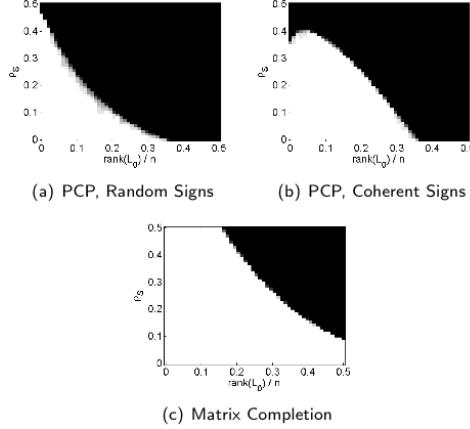


FIGURE 3. Phase Transitions in Probability of Successful Recovery

4. Tyler's M-estimator

In this part we introduce a simple robust covariance estimator, *Tyler's M-estimator*. In general, Huber's M-estimators [Hub81] are generalizations of the MLE estimators to achieve additional properties such as robustness against outliers. M-estimators of covariance are motivated from the MLE estimators under the assumption that data samples are i.i.d. drawn from the elliptical distribution

$x_i \sim C(\rho)e^{-\rho(x^T \Sigma x)} / \sqrt{\det(\Sigma)}$, where $C(\rho)$ is a normalization constant. That is, M-estimator of covariance is defined as the minimizer of

$$L(\Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(x_i^T \Sigma x_i) + \frac{1}{2} \log \det \Sigma$$

By choosing $\rho(t)$ to be some heavy tail distributions, one can achieve robust estimators against outliers.

Tyler's M-estimator [Tyl87] is a special case of the M-estimators of covariance with $\rho(t) = \frac{p}{2} \log(t)$, which allows more possibility for large deviations than normal distribution with $\rho(t) \sim t$. Due to the scale invariance $L(\Sigma) = cL(\Sigma)$ in this case, one often adds the constraint $\text{trace}(\Sigma) = 1$ to make the minimizer unique, i.e.

$$(79) \quad \hat{\Sigma}^{\text{Tyler}} = \arg \min_{\text{trace}(\Sigma)=1, \Sigma \succeq 0} L(\Sigma) := \frac{1}{n} \sum_{i=1}^n \log(x_i^T \Sigma x_i) + \frac{1}{2} \log \det \Sigma$$

The following simple iterative algorithm is also given in [Tyl87],

$$(80) \quad \Sigma^{k+1} = [\text{trace}(S_k)]^{-1} \cdot S_k, \quad \text{where } S_k := \sum_i \frac{x_i x_i^T}{x_i^T [\Sigma^k]^{-1} x_i}.$$

In particular, Tyler showed that it is the “most robust” estimator of the scatter matrix of an elliptical distribution in the sense of minimizing the maximum asymptotic variance. Therefore, we can expect that it leads to robust PCA. For more details, see [Zha16, ZCS14].

5. Sparse PCA

Sparse PCA is firstly proposed by [ZHT06] which tries to locate sparse principal components, which also has a SDP relaxation.

Recall that classical PCA is to solve

$$\begin{aligned} & \max \quad x^T \Sigma x \\ & \text{s.t.} \quad \|x\|_2 = 1 \end{aligned}$$

which gives the maximal variation direction of covariance matrix Σ .

Note that $x^T \Sigma x = \text{trace}(\Sigma(x x^T))$. Classical PCA can thus be written as

$$\begin{aligned} & \max \quad \text{trace}(\Sigma X) \\ & \text{s.t.} \quad \text{trace}(X) = 1 \\ & \quad X \succeq 0 \end{aligned}$$

The optimal solution gives a rank-1 X along the first principal component. A recursive application of the algorithm may lead to top k principal components. That is, one first to find a rank-1 approximation of Σ and extract it from $\Sigma_0 = \Sigma$ to get $\Sigma_1 = \Sigma - X$, then pursue the rank-1 approximation of Σ_1 , and so on.

Now we are looking for sparse principal components, i.e. $\#\{X_{ij} \neq 0\}$ are small. Using 1-norm convexification, we have the following SDP formulation [dGJL07] for Sparse PCA

$$\begin{aligned} & \max \quad \text{trace}(\Sigma X) - \lambda \|X\|_1 \\ & \text{s.t.} \quad \text{trace}(X) = 1 \\ & \quad X \succeq 0 \end{aligned}$$

The following Matlab codes realized the SDP algorithm above by CVX (<http://cvxr.com/cvx>).

```
% Construct a 10-by-20 Gaussian random matrix and form a 20-by-20 correlation
% (inner product) matrix R
X0 = randn(10,20);
R = X0'*X0;

d = 20;
e = ones(d,1);

% Call CVX to solve the SPCA given R
if exist('cvx_setup.m','file'),
    cd /matlab_tools/cvx/
    cvx_setup
end

lambda = 0.5;
k = 10;

cvx_begin
    variable X(d,d) symmetric;
    X == semidefinite(d);
    minimize(-trace(R*X)+lambda*(e'*abs(X)*e));
    subject to
        trace(X)==1;
cvx_end
```

6. MDS with Incomplete Information

In this lecture, we introduce Semi-Definite Programming (SDP) approach to solve some generalized Multi-dimensional Scaling (MDS) problems with uncertainty. Recall that in classical MDS, given pairwise distances $d_{ij} = \|x_i - x_j\|^2$ among a set of points $x_i \in \mathbb{R}^p$ ($i = 1, 2, \dots, n$) whose coordinates are unknown, our purpose is to find $y_i \in \mathbb{R}^k$ ($k \leq p$) such that

$$(81) \quad \min \sum_{i,j=1}^n (\|y_i - y_j\|^2 - d_{ij})^2.$$

In classical MDS (Section 3 in Chapter 1) an eigen-decomposition approach is pursued to find a solution when all pairwise distances d_{ij} 's are known and noise-free. In case that d_{ij} 's are not from pairwise distances, we often use gradient descend method to solve it. However there is no guarantee that gradient descent will converge to the global optimal solution. In this section we will introduce a method based on convex relaxation, in particular the semi-definite relaxation, which will guarantee us to find optimal solutions in the following scenarios.

- Noisy perturbations: $d_{ij} \rightarrow \widetilde{d}_{ij} = d_{ij} + \epsilon_{ij}$

- Incomplete measurements: only partial pairwise distance measurements are available on an edge set of graph, i.e. $G = (V, E)$ and d_{ij} is given when $(i, j) \in E$ (e.g. x_i and x_j in a neighborhood).
- Anchors: sometimes we may fix the locations of some points called *anchors*, e.g. in sensor network localization (SNL) problem.

In other words, we are looking for MDS on graphs with partial and noisy information.

6.1. SD Relaxation of MDS. Like PCA, classical MDS has a semi-definite relaxation. In the following we shall introduce how the constraint

$$(82) \quad \|y_i - y_j\|^2 = d_{ij},$$

can be relaxed into linear matrix inequality system with positive semidefinite variables.

Denote $Y = [y_1, \dots, y_n]^{k \times n}$ where $y_i \in \mathbb{R}^k$, and

$$e_i = (0, 0, \dots, 1, 0, \dots, 0) \in \mathbb{R}^n.$$

Then we have

$$\|y_i - y_j\|^2 = (y_i - y_j)^T (y_i - y_j) = (e_i - e_j)^T Y^T Y (e_i - e_j)$$

Set $X = Y^T Y$, which is symmetric and positive semi-definite. Then

$$\|Y_i - Y_j\|^2 = (e_i - e_j)^T (e_i - e_j) \bullet X.$$

So

$$\|Y_i - Y_j\|^2 = d_{ij}^2 \Leftrightarrow (e_i - e_j)^T (e_i - e_j) \bullet X = d_{ij}^2$$

which is linear with respect to X .

Now we relax the constraint $X = Y^T Y$ to

$$X \succeq Y^T Y \Leftrightarrow X - Y^T Y \succeq 0.$$

Through Schur Complement Lemma we know

$$X - Y^T Y \succeq 0 \Leftrightarrow \begin{bmatrix} I & Y \\ Y^T & X \end{bmatrix} \succeq 0$$

We may define a new variable

$$Z \in S^{k+n}, Z = \begin{bmatrix} I_k & Y \\ Y^T & X \end{bmatrix}$$

which gives the following result.

LEMMA 6.1. The quadratic constraint

$$\|y_i - y_j\|^2 = d_{ij}^2, \quad (i, j) \in E$$

has a semi-definite relaxation:

$$\begin{cases} Z_{1:k, 1:k} = I \\ (0; e_i - e_j)(0; e_i - e_j)^T \bullet Z = d_{ij}^2, \quad (i, j) \in E \\ Z = \begin{bmatrix} I_k & Y \\ Y^T & X \end{bmatrix} \succeq 0. \end{cases}$$

where \bullet denotes the Hadamard inner product, i.e. $A \bullet B := \sum_{i,j=1}^n A_{ij} B_{ij}$.

Note that the constraint with equalities of d_{ij}^2 can be replaced by inequalities such as $\leq d_{ij}^2(1 + \epsilon)$ (or $\geq d_{ij}^2(1 - \epsilon)$). This is a system of linear matrix (in)-equalities with positive semidefinite variable Z . Therefore, the problem becomes a typical semidefinite programming.

Given such a SD relaxation, we can easily generalize classical MDS to the scenarios in the introduction. For example, consider the generalized MDS with anchors which is often called *sensor network localization* problem in literature [BLT⁺06]. Given anchors a_k ($k = 1, \dots, s$) with known coordinates, find x_i such that

- $\|x_i - x_j\|^2 = d_{ij}^2$ where $(i, j) \in E_x$ and x_i are unknown locations
- $\|a_k - x_j\|^2 = \widehat{d}_{kj}^2$ where $(k, j) \in E_a$ and a_k are known locations

We can exploit the following SD relaxation:

- $(0; e_i - e_j)(0; e_i - e_j)^T \bullet Z = d_{ij}$ for $(i, j) \in E_x$,
- $(a_i; e_j)(a_i; e_j)^T \bullet Z = \widehat{d}_{ij}$ for $(i, j) \in E_a$,

both of which are linear with respect to Z .

Recall that every SDP problem has a dual problem (SDD). The SDD associated with the primal problem above is

$$(83) \quad \min \quad I \bullet V + \sum_{i,j \in E_x} w_{ij} d_{ij} + \sum_{i,j \in E_a} \widehat{w}_{ij} \widehat{d}_{ij}$$

s.t.

$$S = \begin{pmatrix} V & 0 \\ 0 & 0 \end{pmatrix} + \sum_{i,j \in E_x} w_{ij} A_{ij} + \sum_{i,j \in E_a} \widehat{w}_{ij} \widehat{A}_{ij} \succeq 0$$

where

$$\begin{aligned} A_{ij} &= (0; e_i - e_j)(0; e_i - e_j)^T \\ \widehat{A}_{ij} &= (a_i; e_j)(a_i; e_j)^T. \end{aligned}$$

The variables w_{ij} is the stress matrix on edge between unknown points i and j and \widehat{w}_{ij} is the stress matrix on edge between anchor i and unknown point j . Note that the dual is always feasible, as $V = 0$, $y_{ij} = 0$ for all $(i, j) \in E_x$ and $w_{ij} = 0$ for all $(i, j) \in E_a$ is a feasible solution.

There are many matlab toolboxes for SDP, *e.g.* CVX, SEDUMI, and recent toolboxes SNLSDP (<http://www.math.nus.edu.sg/~mattohkc/SNLSDP.html>) and DISCO (<http://www.math.nus.edu.sg/~mattohkc/disco.html>) by Toh *et. al.*, adapted to MDS with uncertainty.

A crucial theoretical question is to ask, when $X = Y^T Y$ holds such that SDP embedding Y gives the same answer as the classical MDS? Before looking for answers to this question, we first present an application example of SDP embedding.

6.2. Protein 3D Structure Reconstruction. Here we show an example of using SDP to find 3-D coordinates of a protein molecule based on noisy pairwise distances for atoms in ϵ -neighbors. We use matlab package SNLSDP by Kim-Chuan Toh, Pratik Biswas, and Yinyu Ye, downloadable at <http://www.math.nus.edu.sg/~mattohkc/SNLSDP.html>.

After installation, Figure 4 shows the results of the following codes.

```
>> startup
>> testSNLsolver
```

```
number of anchors = 0
```

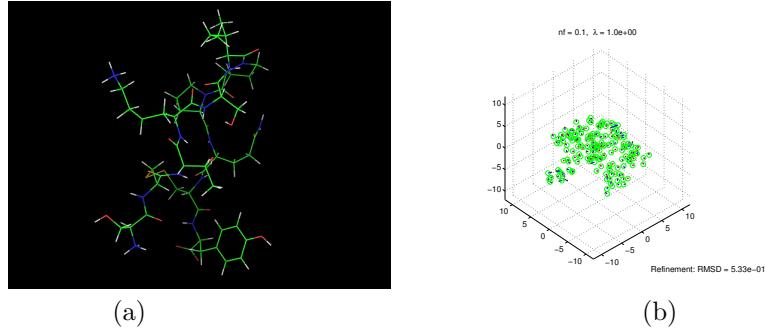


FIGURE 4. (a) 3D Protein structure of PDB-1GM2, edges are chemical bonds between atoms. (b) Recovery of 3D coordinates from SNLSDP with 5 Å-neighbor graph and multiplicative noise at 0.1 level. Red point: estimated position of unknown atom. Green circle: actual position of unknown atom. Blue line: deviation from estimation to the actual position.

```
number of sensors = 166
box scale = 20.00
radius = 5.00
multiplicative noise, noise factor = 1.00e-01
-----
estimate sensor positions by SDP
-----
num of constraints = 2552,
Please wait:
solving SDP by the SDPT3 software package
sdplib = -3.341e+03, time = 34.2s
RMSD = 7.19e-01
-----
refine positions by steepest descent
-----
objstart = 4.2408e+02, objend = 2.7245e+02
number of iterations = 689, time = 0.9s
RMSD = 5.33e-01
-----
(noise factor)^2 = -20.0dB,
mean square error (MSE) in estimated positions = -5.0dB
```

7. Exact Reconstruction and Universal Rigidity

Now we are going to answer the fundamental question, when the SDP relaxation exactly reconstruct the coordinates up to a rigid transformation. We will provide two theories, one from the optimality rank properties of SDP, and the other from a geometric criterion, *universal rigidity*.

Recall that for a standard SDP with $X, C \in \mathbb{R}^{n \times n}$

$$(84) \quad \begin{aligned} \min \quad & C \bullet X = \sum_{i,j} c_{ij} X_{ij} \\ \text{s.t.} \quad & A_i \bullet X = b_i, \quad \text{for } i = 1, \dots, m \\ & X \succeq 0 \end{aligned}$$

whose SDD is

$$(85) \quad \begin{aligned} \max \quad & b^T y \\ \text{s.t.} \quad & S = C - \sum_{i=1}^m A_i y_i \succeq 0. \end{aligned}$$

Such SDP has the following rank properties [Ali95]:

- A. maximal rank solutions X^* or S^* exist;
- B. minimal rank solutions X^* or S^* exist;
- C. if complementary condition $X^* S^* = 0$ holds, then $\text{rank}(X^*) + \text{rank}(S^*) \leq n$ with equality holds iff strictly complementary condition holds, whence $\text{rank}(S^*) \geq n - k \Rightarrow \text{rank}(X^*) \leq k$.

Strong duality of SDP tells us that an interior point feasible solution in primal or dual problem will ensure the complementary condition and the zero duality gap. Now we assume that $d_{ij} = \|x_i - x_j\|$ precisely for some unknown $x_i \in \mathbb{R}^k$. Then the primal problem is feasible with $Z = (I_d; Y)^T (I_d; Y)$. Therefore the complementary condition holds and the duality gap is zero. In this case, assume that Z^* is a primal feasible solution of SDP embedding and S^* is an optimal dual solution, then

- (1) $\text{rank}(Z^*) + \text{rank}(S^*) \leq k + n$ and $\text{rank}(Z^*) \geq k$, whence $\text{rank}(S^*) \leq n$;
- (2) $\text{rank}(Z^*) = k \iff X = Y^T Y$.

It follows that if an optimal dual S^* has rank n , then every primal solution Z^* has rank k , which ensures $X = Y^T Y$. Therefore it suffices to find a maximal rank dual solution S^* whose rank is n .

Above we have optimality rank condition from SDP. Now we introduce a geometric criterion based on universal rigidity.

DEFINITION (Universal Rigidity (UR) or Unique Localization (UL)). $\exists! y_i \in \mathbb{R}^k \hookrightarrow \mathbb{R}^l$ where $l \geq k$ s.t. $d_{ij}^2 = \|y_i - y_j\|^2, \widehat{d}_{ij}^2 = \|a_k - y_j\|^2$.

It simply says that there is no nontrivial extension of $y_i \in \mathbb{R}^k$ in \mathbb{R}^l satisfying $d_{ij}^2 = \|y_i - y_j\|^2$ and $\widehat{d}_{ij}^2 = \|(a_k; 0) - y_j\|^2$. The following is a short history about universal rigidity.

[Schoenberg 1938] G is complete \implies UR

[So-Ye 2007] G is incomplete \implies UR \iff SDP has maximal rank solution $\text{rank}(Z^*) = k$.

THEOREM 7.1. [SY07] The following statements are equivalent.

- (1) The graph is universally rigid or has a unique localization in \mathbb{R}^k .
- (2) The max-rank feasible solution of the SDP relaxation has rank k ;
- (3) The solution matrix has $X = Y^T Y$ or $\text{trace}(X - Y^T Y) = 0$.

Moreover, the localization of a UR instance can be computed approximately in a time polynomial in n, k , and the accuracy $\log(1/\epsilon)$.

In fact, the max-rank solution of SDP embedding is unique. There are many open problems in characterizing UR conditions, see Ye's survey at ICCM'2010.

In practice, we often meet problems with noisy measurements $\alpha d_{ij}^2 \geq \tilde{d}_{ij}^2 \leq \beta d_{ij}^2$. If we relax the constraint $\|y_i - y_j\|^2 = d_{ij}^2$ or equivalently $A_i \bullet X = b_i$ to inequalities, however we can achieve arbitrary small rank solution. To see this, assume that

$$A_i X = b_i \quad \mapsto \quad \alpha b_i \leq A_i X \leq \beta b_i \quad i = 1, \dots, m, \text{ where } \beta \geq 1, \alpha \in (0, 1)$$

then So, Ye, and Zhang (2008) [SYZ08] show the following result.

THEOREM 7.2. For every $d \geq 1$, there is a SDP solution $\hat{X} \succeq 0$ with rank $\text{rank}(\hat{X}) \leq d$, if the following holds,

$$\beta = \begin{cases} 1 + \frac{18 \ln 2m}{d} & 1 \leq d \leq 18 \ln 2m \\ 1 + \frac{\sqrt{18 \ln 2m}}{d} & d \geq 18 \ln 2m \end{cases}$$

$$\alpha = \begin{cases} \frac{1}{e(2m)^{2/d}} & 1 \leq d \leq 4 \ln 2m \\ \max \left\{ \frac{1}{e(2m)^{2/d}}, 1 - \sqrt{\frac{4 \ln 2m}{d}} \right\} & d \geq 4 \ln 2m \end{cases}$$

Note that α, β are independent to n .

8. Maximal Variance Unfolding

Here we give a special case of SDP embedding, Maximal Variance Unfolding (MVU) [WS06]. In this case we choose graph $G = (V, E)$ as k -nearest neighbor graph. As a contrast to the SDP embedding above, we did not pursue a semi-definite relaxation $X \succeq Y^T Y$, but instead define it as a positive semi-definite kernel $K = Y^T Y$ and maximize the trace of K .

Consider a set of points x_i ($i = 1, \dots, n$) whose pairwise distance d_{ij} is known if x_j lies in k -nearest neighbors of x_i . In other words, consider a k -nearest neighbor graph $G = (V, E)$ with $V = \{x_i : i = 1, \dots, n\}$ and $(i, j) \in E$ if j is a member of k -nearest neighbors of i .

Our purpose is to find coordinates $y_i \in \mathbb{R}^k$ for $i = 1, 2, \dots, n$ s.t.

$$d_{ij}^2 = \|y_i - y_j\|^2$$

wherever $(i, j) \in E$ and $\sum_i y_i = 0$.

Set $K_{ij} = \langle y_i, y_j \rangle$. Then K is symmetric and positive semidefinite, which satisfies

$$K_{ii} + K_{jj} - 2K_{ij} = d_{ij}^2.$$

There are possibly many solutions for such K , and we look for a particular one with maximal trace which characterizes the maximal variance.

$$(86) \quad \begin{aligned} \max \quad & \text{trace}(K) = \sum_{i=1}^n \lambda_i(K) \\ \text{s.t.} \quad & K_{ii} + K_{jj} - 2K_{ij} = d_{ij}^2, \\ & \sum_j K_{ij} = 0, \\ & K \succeq 0 \end{aligned}$$

Again it is a SDP. The final embedding is obtained by using eigenvector decomposition of $K = Y^T Y$.

However we note here that maximization of trace is not a provably good approach to “unfold” a manifold. Sometimes, there are better ways than MVU, *e.g.* if original data lie on a plane then maximization of the diagonal distance between two neighboring triangles will unfold and force it to be a plane. This is a special case of the general $k+1$ -lateration graphs [SY07]. From here we see that there are other linear objective functions better than trace for the purpose of “unfolding” a manifold.

CHAPTER 5

Manifold Learning

1. Introduction

In the past month we talked about two topics: one is the sample mean and sample covariance matrix (PCA) in high dimensional spaces. We have learned that when dimension p is large and sample size n is relatively small, in contrast to the traditional statistics where p is fixed and $n \rightarrow \infty$, both sample mean and PCA may have problems. In particular, Stein's phenomenon shows that in high dimensional space with independent Gaussian distributions, the sample mean is worse than a shrinkage estimator; moreover, random matrix theory sheds light on that in high dimensional space with sample size in a fixed ratio of dimension, the sample covariance matrix and PCA may not reflect the signal faithfully. These phenomena start a new philosophy in high dimensional data analysis that to overcome the curse of dimensionality, additional constraints has to be put that data never distribute in every corner in high dimensional spaces. Sparsity is a common assumption in modern high dimensional statistics. For example, data variation may only depend on a small number of variables; independence of Gaussian random fields leads to sparse covariance matrix; and the assumption of conditional independence can also lead to sparse inverse covariance matrix. In particular, an assumption that data concentrate around a low dimensional manifold in high dimensional spaces, leads to manifold learning or nonlinear dimensionality reduction, e.g. ISOMAP, LLE, and Diffusion Maps etc. This assumption often finds example in computer vision, graphics, and image processing.

All the work introduced in this chapter can be regarded as generalized PCA/MDS on nearest neighbor graphs, which has roots in manifold learning concept. Two pieces of milestone works, ISOMAP [TdSL00] and Locally Linear Embedding (LLE) [RL00], are firstly published in *science* 2000, which opens a new field called nonlinear dimensionality reduction, or manifold learning in high dimensional data analysis. Here is the development of manifold learning method:

$$(87) \quad \text{MDS} \longrightarrow \text{ISOMAP}$$

$$\text{PCA} \longrightarrow \text{LLE} \longrightarrow \left\{ \begin{array}{l} \text{Local Tangent Space Alignment} \\ \text{Hessian LLE} \\ \text{Laplacian Eigen Map} \\ \text{Diffusion Map} \end{array} \right.$$

To understand the motivation of such a novel methodology, let's take a brief review on PCA/MDS. Given a set of data $x_i \in \mathbb{R}^p$ ($i = 1, \dots, n$) or merely pairwise distances $d(x_i, x_j)$, PCA/MDS essentially looks for an affine space which best capture the variation of data distribution, see Figure 1(a). However, this scheme will not work in the scenario that data are actually distributed on a highly nonlinear

curved surface, i.e. *manifolds*, see the example of Swiss Roll in Figure 1(b). Can we extend PCA/MDS in certain sense to capture intrinsic coordinate systems which charts the manifold?

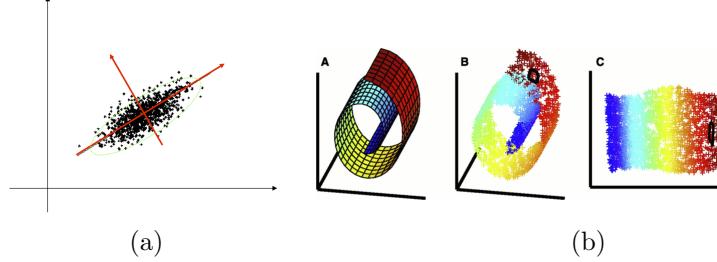


FIGURE 1. (a) Find an affine space to approximate data variation in PCA/MDS. (b) Swiss Roll data distributed on a nonlinear 2-D submanifold in Euclidean space \mathbb{R}^3 . Our purpose is to capture an intrinsic coordinate system describing the submanifold.

ISOMAP and LLE, as extensions from MDS and local PCA, respectively, leads to a series of attempts to address this problem.

All the current techniques in manifold learning, as extensions of PCA and MDS, are often called as *Spectral Kernel Embedding*. The common theme of these techniques can be described in Figure 2. The basic problem is: given a set of data points $\{x_1, x_2, \dots, x_n \in \mathbb{R}^p\}$, how to find out $y_1, y_2, \dots, y_n \in \mathbb{R}^d$, where $d \ll p$, such that some geometric structures (local or global) among data points are best preserved.

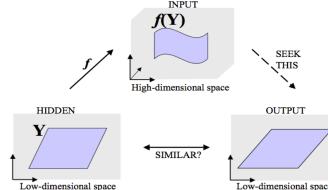


FIGURE 2. The generative model for manifold learning. Y is the hidden parameter space (like rotation angle of faces below), f is a measure process which maps Y into a sub-manifold in a high dimensional ambient space, $X = f(Y) \subset \mathbb{R}^p$. All of our purpose is to recover this hidden parameter space Y given samples $\{x_i \in \mathbb{R}^p : i = 1, \dots, n\}$.

All the manifold learning techniques can be summarized in the following meta-algorithm, which explains precisely the name of *spectral kernel embedding*. All the methods can be called certain *eigenmaps* associated with some positive semi-definite kernels.

- (1) Construct a data graph $G = (V, E)$, where $V = \{x_i : i = 1, \dots, n\}$. For example,

ε -neighborhood, $i \sim j \Leftrightarrow d(x_i, x_j) \leq \varepsilon$, which leads to an undirected graph;

k -nearest neighbor, $(i, j) \in E \Leftrightarrow j \in \mathcal{N}_k(i)$, which leads to a directed graph.

- (2) Construct a positive semi-definite matrix K (kernel).
- (3) Eigen-decomposition $K = U \Lambda U^T$, then $Y_d = U_d \Lambda_d^{\frac{1}{2}}$, where choose d eigenvectors (top or bottom) U_d .

EXAMPLE 7 (PCA). G is complete, $K = \hat{\Sigma}_n$ is a covariance matrix.

EXAMPLE 8 (MDS). G is complete, $K = -\frac{1}{2} HDH^T$, where $D_{ij} = d^2(x_i, x_j)$.

EXAMPLE 9 (ISOMAP). G is incomplete.

$$D_{ij} = \begin{cases} d(x_i, x_j) & \text{if } (i, j) \in E, \\ \hat{d}_g(x_i, x_j) & \text{if } (i, j) \notin E. \end{cases}$$

where \hat{d}_g is a graph shortest path. Then

$$K = -\frac{1}{2} HDH^T.$$

Note that K is positive semi-definite if and only if D is a squared distance matrix.

EXAMPLE 10 (LLE). G is incomplete. $K = (I - W)^T(I - W)$, where

$$W_{ij}^{n \times n} = \begin{cases} w_{ij} & j \in \mathcal{N}(i), \\ 0 & \text{other's.} \end{cases}$$

and w_{ij} solves the following optimization problem

$$\min_{\sum_j w_{ij}=1} \|X_i - \sum_{j \in \mathcal{N}(i)} w_{ij} X_j\|^2.$$

After obtaining W , compute the global embedding d -by- n embedding matrix $Y = [Y_1, \dots, Y_n]$,

$$\min_Y \sum_{i=1}^n \|Y_i - \sum_{j=1}^n W_{ij} Y_j\|^2 = \text{trace}((I - W)Y^T Y(I - W)^T).$$

This is equivalent to find smallest eigenvectors of $K = (I - W)^T(I - W)$.

2. ISOMAP

ISOMAP is an extension of MDS, where pairwise euclidean distances between data points are replaced by geodesic distances, computed by *graph shortest path distances*.

- (1) Construct a neighborhood graph $G = (V, E, d_{ij})$ such that
 - $V = \{x_i : i = 1, \dots, n\}$
 - $E = \{(i, j) : \text{if } j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors, ϵ -neighbors
 - $d_{ij} = d(x_i, x_j)$, e.g. Euclidean distance when $x_i \in \mathbb{R}^p$
- (2) Compute graph shortest path distances
 - $d_{ij} = \min_{P=(x_i, \dots, x_j)} (\|x_i - x_{t_1}\| + \dots + \|x_{t_{k-1}} - x_j\|)$, is the length of a graph shortest path connecting i and j
 - Dijkstra's algorithm ($O(kn^2 \log n)$) and Floyd's Algorithm ($O(n^3)$)

(3) classical MDS with $D = (d_{ij}^2)$

construct a symmetric (positive semi-definite if D is a squared distance) $B = -0.5H D H^T$ where $H = I - \mathbf{1}\mathbf{1}^T/n$ (or $H = I - \mathbf{1}\mathbf{a}^T$ for any $\mathbf{a}^T\mathbf{1} = 1$).

Find eigenvector decomposition of $B = U\Lambda U^T$ and choose top d eigenvectors as embedding coordinates in \mathbb{R}^d , i.e. $Y_d = [y_1, \dots, y_d] = [U_1, \dots, U_d]\Lambda_d^{1/2} \in \mathbb{R}^{n \times d}$

Algorithm 6: ISOMAP Algorithm

Input: Metric distance $d_{ij} = d(x_i, x_j)$ between data points and a weighted graph $G = (V, E)$ such that

1 $V = \{x_i : i = 1, \dots, n\}$

2 $E = \{(i, j) : \text{if } j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors, ϵ -neighbors

3 $d_{ij} = d(x_i, x_j)$ for $(i, j) \in E$, e.g. Euclidean distance when $x_i \in \mathbb{R}^p$

Output: Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{k \times n}$ of data.

4 **Step 1:** Compute graph shortest path distances

$$d_{ij} = \min_{P=(x_i, \dots, x_j)} (\|x_i - x_{t_1}\| + \dots + \|x_{t_{k-1}} - x_j\|),$$

which is the length of a graph shortest path connecting i and j ;

5 **Step 2:** Compute $K = -\frac{1}{2}H \cdot D \cdot H^T$ ($D := [d_{ij}^2]$), where H is the Householder centering matrix;

6 **Step 3:** Compute Eigenvalue decomposition $K = U\Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$;

7 **Step 4:** Choose top d nonzero eigenvalues and corresponding eigenvectors,

$$Y_d = U_d \Lambda_d^{1/2} \text{ where}$$

$$U_d = [u_1, \dots, u_d], \quad u_j \in \mathbb{R}^n,$$

$$\Lambda_d = \text{diag}(\lambda_1, \dots, \lambda_d)$$

with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$.

The basic feature of ISOMAP can be described as: we find a low dimensional embedding of data such that points nearby are mapped nearby and points far away are mapped far away. In other words, we have global control on the data distance and the method is thus a *global* method. The major shortcoming of ISOMAP lies in its computational complexity, characterized by a full matrix eigenvector decomposition.

2.1. ISOMAP Example. Now we give an example of ISOMAP with matlab codes.

```
% load 33-face data
load ../data/face.mat Y
X = reshape(Y, [size(Y,1)*size(Y,2) size(Y,3)]);
p = size(X,1);
n = size(X,2);
D = pdist(X');
DD = squareform(D);

% ISOMAP embedding with 5-nearest neighbors
```

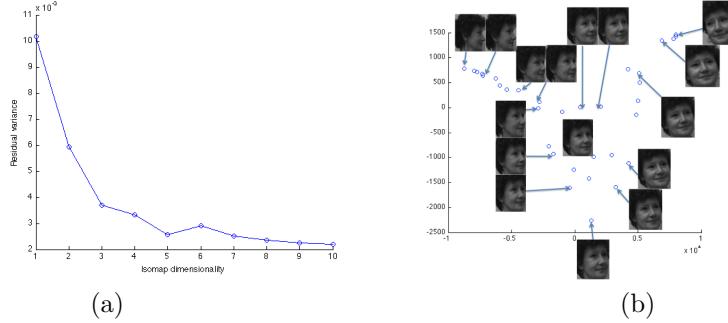


FIGURE 3. (a) Residual Variance plot for ISOMAP. (b) 2-D ISOMAP embedding, where the first coordinate follows the order of rotation angles of the face.

```
[Y_iso,R_iso,E_iso]=isomapII(DD,'k',5);
```

```
% Scatter plot of top 2-D embeddings
y=Y_iso.coords{2};
scatter(y(1,:),y(2,:))
```

2.2. Convergence of ISOMAP. Under dense-sample and regularity conditions on manifolds, ISOMAP is proved to show convergence to preserve geodesic distances on manifolds. The key is to approximate geodesic distance on manifold by a sequence of short Euclidean distance hops.

Consider arbitrary two points on manifold $x, y \in M$. Define

$$\begin{aligned} d_M(x, y) &= \inf_{\gamma} \{\text{length}(\gamma)\} \\ d_G(x, y) &= \min_P (\|x_0 - x_1\| + \dots + \|x_{t-1} - x_t\|) \\ d_S(x, y) &= \min_P (d_M(x_0, x_1) + \dots + d_M(x_{t-1}, x_t)) \end{aligned}$$

where γ varies over the set of smooth arcs connecting x to y in M and P varies over all paths along the edges of G starting at $x_0 = x$ and ending at $x_t = y$. We are going to show $d_M \approx d_G$ with the bridge d_S .

It is easy to see the following upper bounds by d_S :

$$(88) \quad d_M(x, y) \leq d_S(x, y)$$

$$(89) \quad d_G(x, y) \leq d_S(x, y)$$

where the first upper bound is due to triangle inequality for the metric d_M and the second upper bound is due to that Euclidean distances $\|x_i - x_{i+1}\|$ are smaller than arc-length $d_M(x_i, x_{i+1})$.

To see other directions, one has to impose additional conditions on sample density and regularity of manifolds.

LEMMA 2.1 (Sufficient Sampling). Let $G = (V, E)$ where $V = \{x_i : i = 1, \dots, n\} \subseteq M$ is a ϵ -net of manifold M , i.e. for every $x \in M$ there exists $x_i \in V$

such that $d_M(x, x_i) < \epsilon$, and $\{i, j\} \in E$ if $d_M(x_i, x_j) \leq \alpha\epsilon$ ($\alpha \geq 4$). Then for any pair $x, y \in V$,

$$d_S(x, y) \leq \max(\alpha - 1, \frac{\alpha}{\alpha - 2})d_M(x, y).$$

PROOF. Let γ be a shortest path connecting x and y on M whose length is l . If $l \leq (\alpha - 2)\epsilon$, then there is an edge connecting x and y whence $d_S(x, y) = d_M(x, y)$. Otherwise split γ into pieces such that $l = l_0 + tl_1$ where $l_1 = (\alpha - 2)\epsilon$ and $\epsilon \leq l_0 < (\alpha - 2)\epsilon$. This divides arc γ into a sequence of points $\gamma_0 = x, \gamma_1, \dots, \gamma_{t+1} = y$ such that $d_M(x, \gamma_1) = l_0$ and $d_M(\gamma_i, \gamma_{i+1}) = l_1$ ($i \geq 1$). There exists a sequence of $x_0 = x, x_1, \dots, x_{t+1} = y$ such that $d_M(x_i, \gamma_i) \leq \epsilon$ and

$$\begin{aligned} d_M(x_i, x_{i+1}) &\leq d_M(x_i, \gamma_i) + d_M(\gamma_i, \gamma_{i+1}) + d_M(\gamma_{i+1}, x_{i+1}) \\ &\leq \epsilon + l_1 + \epsilon \\ &= \alpha\epsilon \\ &= l_1\alpha/(\alpha - 2) \end{aligned}$$

whence $(x_i, x_{i+1}) \in E$. Similarly $d_M(x, x_1) \leq d_M(x, \gamma_1) + d_M(\gamma_1, x_1) \leq (\alpha - 1)\epsilon \leq l_0(\alpha - 1)$.

$$\begin{aligned} d_S(x, y) &\leq \sum_{i=0}^{t-1} d_M(x_i, x_{i+1}) \\ &\leq l \max\left(\frac{\alpha}{\alpha - 2}, \alpha - 1\right) \end{aligned}$$

Setting $\alpha = 4$ gives rise to $d_S(x, y) \leq 3d_M(x, y)$. □

The other lower bound $d_S(x, y) \leq cd_G(x, y)$ requires that for every two points x_i and x_j , Euclidean distance $\|x_i - x_j\| \leq cd_M(x_i, x_j)$. This imposes a regularity on manifold M , whose curvature has to be bounded. We omit this part here and leave the interested readers to the reference by Bernstein, de Silva, Langford, and Tenenbaum 2000, as a supporting information to the ISOMAP paper.

3. Locally Linear Embedding (LLE)

In applications points nearby should be mapped nearby, while points far away should impose no constraint. This is because typically when points are close enough, they are similar, while points are far, there is no faithful information to measure how far they are. Therefore global information about geodesic distance might not be accurate, in addition to its expensive computational cost. This motivates another type of algorithm, locally linear embedding. The algorithm assume that any data point in a high dimensional ambient space can be a linear combination of data points in its neighborhood. In other words, a data point x_i has its neighborhood $x_j \in \mathcal{N}_i$ deciding its sufficient statistics. Alignment of such local linear structures can lead to a global unfolding of data manifolds, often described as *fit locally and think globally*. This is a *local* method as it involves data points in local neighbors and hence a sparse eigenvector decomposition.

Now we are going to describe the procedure of LLE.

The reason behind the crucial steps can be explained as follows.

- (1) Local fitting:

Algorithm 7: LLE Algorithm

-
- Input:** A graph $G = (V, E)$ such that
- 1 $V = \{x_i : i = 1, \dots, n\}$
 - 2 $E = \{(i, j) : j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors, ϵ -neighbors
 - Output:** Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{d \times n}$ of data.
 - 3 **Step 1** (local fitting): for each x_i and its neighbors \mathcal{N}_i , solve

4

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij} = 1} \|x_i - \sum_{j \in \mathcal{N}_i} w_{ij} x_j\|^2,$$

by $\hat{w}_i(\mu) = (C_i + \mu I)^{-1} \mathbf{1}$ for some regularization parameter $\mu > 0$ and
 $w_i = \hat{w}_i / \hat{w}_i^T \mathbf{1}$;

- 5 **Step 2** (global alignment): define the weight embedding matrix

6

$$W_{ij} = \begin{cases} w_{ij}, & j \in \mathcal{N}_i \\ 0, & \text{otherwise} \end{cases}$$

Compute $K = (I - W)^T (I - W)$ which is a positive semi-definite kernel matrix;

- 7 **Step 3** (Eigenmap): Compute Eigenvalue decomposition $K = U \Lambda U^T$ with
 $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \lambda_{n-1} > \lambda_n = 0$; choose bottom $d+1$ nonzero eigenvalues and corresponding eigenvectors and drop the smallest eigenvalue-eigenvector (0-constant) pair, such that

$$U_d = [u_{n-d}, \dots, u_{n-1}], \quad u_j \in \mathbb{R}^n,$$

$$\Lambda_d = \text{diag}(\lambda_{n-d}, \dots, \lambda_{n-1}).$$

Define $Y_d = U_d \Lambda_d^{\frac{1}{2}}$.

Pick up a point x_i and its neighbors \mathcal{N}_i . Compute the local fitting weights

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij} = 1} \|x_i - \sum_{j \in \mathcal{N}_i} w_{ij} x_j\|^2,$$

which is equivalent to

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij} = 1} \left\| \sum_{j \in \mathcal{N}_i} w_{ij} (x_j - x_i) \right\|^2,$$

that is, finding a linear combination (possibly *not unique!*) for the subspace spanned by $\{(x_j - x_i) : j \in \mathcal{N}_i\}$. This can be done by Lagrange multiplier method, *i.e.* solving

$$\min_{w_{ij}} \frac{1}{2} \left\| \sum_{j \in \mathcal{N}_i} w_{ij} (x_j - x_i) \right\|^2 + \lambda \left(1 - \sum_{j \in \mathcal{N}_i} w_{ij} \right).$$

Let $w_i = [w_{ij_1}, \dots, w_{ij_k}]^T \in \mathbb{R}^k$, $\bar{X}_i = [x_{j_1} - x_i, \dots, x_{j_k} - x_i]$, and the local Gram (covariance) matrix $C_i(j, k) = \langle x_j - x_i, x_k - x_i \rangle$, whence the weights are

$$(90) \quad w_i = \lambda C_i^\dagger \mathbf{1},$$

where the Lagrange multiplier equals to the following normalization parameter

$$(91) \quad \lambda = \frac{1}{\mathbf{1}^T C_i^\dagger \mathbf{1}},$$

TABLE 1. Comparisons between ISOMAP and LLE.

ISOMAP	LLE
MDS on geodesic distance matrix global approach	local PCA + eigen-decomposition local approach
no for nonconvex manifolds with holes landmark (Nystrom)	ok with nonconvex manifolds with holes Hessian
Extensions: conformal isometric, etc.	Extensions: Laplacian LTSA etc.

and C_i^\dagger is a Moore-Penrose (pseudo) inverse of C_i . Note that C_i is often ill-conditioned and to find its Moore-Penrose inverse one can use regularization method $(C_i + \mu I)^{-1}$ for some $\mu > 0$.

(2) Global alignment

Define a n -by- n weight matrix W :

$$W_{ij} = \begin{cases} w_{ij}, & j \in \mathcal{N}_i \\ 0, & \text{otherwise} \end{cases}$$

Compute the global embedding d -by- n embedding matrix Y ,

$$\min_Y \sum_i \|y_i - \sum_{j=1}^n W_{ij} y_j\|^2 = \text{trace}(Y(I - W)^T (I - W) Y^T)$$

In other words, construct a positive semi-definite matrix $K = (I - W)^T (I - W)$ and find $d+1$ smallest eigenvectors of K , v_0, v_1, \dots, v_d associated smallest eigenvalues $\lambda_0, \dots, \lambda_d$. Drop the smallest eigenvector which is the constant vector explaining the degree of freedom as translation and set $Y = [v_1/\sqrt{\lambda_1}, \dots, v_d/\sqrt{\lambda_d}]^T$.

The benefits of LLE are:

- Neighbor graph: k -nearest neighbors is of $O(kn)$
- W is sparse: $kn/n^2 = k/n$ non-zeroes
- $K = (I - W)^T (I - W)$ is guaranteed to be positive semi-definite

However, unlike ISOMAP, it is not clear if LLE constructed above converges under certain conditions. This has to be left to some variations of basic LLE above, such as Laplacian LLE, Hessian LLE, and LTSA etc. with convergence guarantees.

3.1. Issues of LLE and a Modified Version. Using the regularization, (90) leads to a family of weight vectors

$$(92) \quad w_i(\mu) = \lambda(C_i + \mu I)^{-1} \mathbf{1} = \sum_j \frac{1}{\lambda_j^{(i)} + \mu} v_j v_j^T \mathbf{1}$$

where the local PCA $C_i = V \Lambda V^T$ ($\Lambda = \text{diag}(\lambda_j^{(i)})$, $V = [v_j]$).

So basically $w_i(\mu)$ is made up of a low-pass filter: the projections of $\mathbf{1}$ on those directions U_j such that $\lambda_j^{(i)} \ll \mu$ are preserved while those projections with $\lambda_j^{(i)} \gg \mu$ are attenuated. In the ideal case without noise, such a low-pass filter makes $w_i(\mu)$ spanned by the normal subspace orthogonal to the local PCA, such that the reconstructed Y will follow the directions of local PCA. However, in applications when noise are presented, especially not well separated with signals, such $w_i(\mu)$ is

Algorithm 8: MLLE Algorithm

Input: A graph $G = (V, E)$ such that

- 1 $V = \{x_i : i = 1, \dots, n\}$
- 2 $E = \{(i, j) : j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors, ϵ -neighbors

Output: Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{d \times n}$ of data.

- 3 **Step 1** (local fitting): for each x_i and its neighbors \mathcal{N}_i , solve

4

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij}=1} \|x_i - \sum_{j \in \mathcal{N}_i} w_{ij} x_j\|^2,$$

by $\hat{w}_i(\mu) = (C_i + \mu I)^{-1} \mathbf{1}$ for some regularization parameter $\mu > 0$ and

$$w_i = \hat{w}_i / \hat{w}_i^T \mathbf{1};$$

- 5 **Step 2** (local residue PCA): for each x_i and its neighbors \mathcal{N}_i ($k_i = |\mathcal{N}_i|$), let $C_i = V \Lambda V^T$ be its eigenvalue decomposition where $\Lambda = (\lambda_1, \dots, \lambda_{k_i})$ with $\lambda_1 \geq \dots \geq \lambda_{k_i}$. Find the size of almost normal subspace s_i as the maximal size that the ratio of residue eigenvalue sum over principle eigenvalue sum is below a threshold, i.e.

$$s_i = \max_l \left\{ l \leq k_i - d, \frac{\sum_{j=k_i-l+1}^{k_i} \lambda_j}{\sum_{j=1}^{k_i-l} \lambda_j} \leq \eta \right\}$$

where η is a parameter, such as the median of ratios of residue eigenvalue sum over principle eigenvalue sum. Construct the normal subspace basis matrix as s_i -bottom eigenvector matrix of C_i , $V_i = [v_{k_i-s_i+1}, \dots, v_{k_i}] \in \mathbb{R}^{k_i \times s_i}$, define the weight matrix

$$W_i = (1 - \alpha_i) w_i(\mu) \mathbf{1}_{s_i}^T + V_i H_i^T \in \mathbb{R}^{k_i \times s_i},$$

where $\alpha_i = \|V_i^T \mathbf{1}_{k_i}\|_2 / \sqrt{s_i}$ and $H_i = I_{s_i} - 2uu^T / \|u\|^2$ with $u = V_i^T \mathbf{1}_{k_i} - \alpha_i \mathbf{1}_{s_i}$ (or $u = 0$ if it is small).

- 6 **Step 3** (global alignment): define the weight embedding matrix

7

$$\widehat{W}_i(j, :) = \begin{cases} -\mathbf{1}_{s_i}^T, & j = i, \\ W_i, & j \in \mathcal{N}_i, \\ 0, & \text{otherwise.} \end{cases}$$

Compute $K = \widehat{W}^T \widehat{W}$ which is a positive semi-definite kernel matrix;

- 8 **Step 4** (Eigenmap): Compute Eigenvalue decomposition $K = U \Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{n-1} > \lambda_n = 0$; choose bottom $d+1$ nonzero eigenvalues and corresponding eigenvectors and drop the smallest eigenvalue-eigenvector (0-constant) pair, such that

$$U_d = [u_{n-d}, \dots, u_{n-1}], \quad u_j \in \mathbb{R}^n,$$

$$\Lambda_d = \text{diag}(\lambda_{n-d}, \dots, \lambda_{n-1}).$$

Define $Y_d = U_d \Lambda_d^{\frac{1}{2}}$.

sensitive to the noise direction and might be mixed with signal directions. LLE in this case might not capture well the signal directions in local PCA. Hessian LLE and LTSA are both improvements over this by exploiting all the local principal components. On the other hand, Modified Locally Linear Embedding (MLLE) [ZW] remedies the issue using multiple weight vectors projected from orthogonal complement of local PCA.

MLLE replace the weight vector above by a weight matrix $W_i \in \mathbb{R}^{k_i \times s_i}$, a family of s_i weight vectors using bottom s_i eigenvectors of C_i , $V_i = [v_{k_i-s_i+1}, \dots, v_{k_i}] \in$

$\mathbb{R}^{k_i \times s_i}$, such that

$$(93) \quad W_i = (1 - \alpha_i)w_i(\mu)\mathbf{1}_{s_i}^T + V_i H_i^T,$$

where $\alpha_i = \|V_i^T \mathbf{1}_{k_i}\|_2 / \sqrt{s_i}$ and $H_i = I_{s_i} - 2uu^T$ ($\|u\|_2 = 1$ or 0) is a Householder matrix ($H_i := I_{s_i}$ if $u = 0$) such that $HV_i^T \mathbf{1}_{k_i} = \alpha_i \mathbf{1}_{s_i}$ (hence $W_i^T \mathbf{1}_{k_i} = \mathbf{1}_{s_i}$, every column of W_i being a legal weight vector). In fact, one can choose u in the direction of $V_i^T \mathbf{1}_{k_i} - \alpha_i \mathbf{1}_{s_i}$. An adaptive choice of s_i is given in [ZW] using the trade-off between residual variation and explained variation. Equipped with this weight matrix, one can set the objective function by simultaneously minimizing the residue over all reconstruction weights:

$$\min_Y \sum_i \sum_{l=1}^{s_i} \|y_i - \sum_{j \in \mathcal{N}_i} W_i(j, l)y_j\|^2 := \sum_i \|Y \widehat{W}_i\|_F^2 = \text{trace}[Y(\sum_i \widehat{W}_i \widehat{W}_i^T)Y^T]$$

where \widehat{W}_i is the embedding of $W_i \in \mathbb{R}^{k_i \times s_i}$ into $\mathbb{R}^{n \times s_i}$,

$$\widehat{W}_i(j, :) = \begin{cases} -\mathbf{1}_{s_i}^T, & j = i, \\ W_i, & j \in \mathcal{N}_i, \\ 0, & \text{otherwise.} \end{cases}$$

Python `scikit-learn` package contains an implementation of MLLE. The error analysis of MLLE is similar to that of LTSA [ZW], hence it is expected both lead to similar results in applications. Yet due to the adaptive choice of s_i , MLLE can be adaptive to the heterogeneity in manifold curvature variations.

A shortcoming of MLLE lies in its projection of vector $\mathbf{1}$ on to the local normal subspace spanned by the bottom eigenvectors which might be totally contaminated by noise. So such a computation is to capture noise instead of signal and might be very expensive (due to the full spectrum of C_i) when the intrinsic dimensionality is low. On the other hand both Hessian LLE and LTSA only exploit the partial local SVD which is more robust to noise and cheaper in computational cost.

4. Hessian LLE

Hessian LLE only exploits top eigenvectors of local SVD which is more robust to noise than MLLE, and is provable to find a linear coordinate chart for local isometric and nonconvex assumptions while ISOMAP requires global isometry and convexity assumptions.

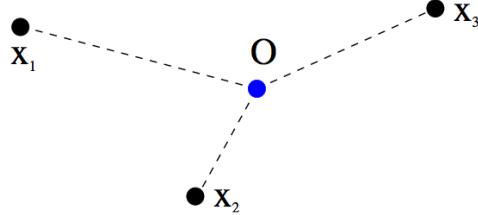


FIGURE 4. Local coordinate system at the origin $O = x_i$.

In LLE, one chooses the weights w_{ij} to minimize the following energy

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij}=1} \left\| \sum_{j \in \mathcal{N}_i} w_{ij}(x_j - x_i) \right\|^2.$$

In the ideal case, if the points $\tilde{x}_j = x_j - x_i$ are linearly dependent, then there is some w_{ij} , possibly not unique, such that $0 = \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_j$. In this local chart (Figure 4), we have

$$0 = \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_j, \quad \text{and} \quad 1 = \sum_{j \in \mathcal{N}_i} w_{ij}.$$

For any smooth function $y(x)$, consider its Taylor expansion up to the second order

$$y(x) = y(0) + x^T \nabla y(0) + \frac{1}{2} x^T (\mathcal{H}y)(0) x + o(\|x\|^2).$$

Therefore

$$\begin{aligned} (I - W)y(0) &:= y(0) - \sum_{j \in \mathcal{N}_i} w_{ij} y(\tilde{x}_j) \\ &\approx y(0) - \sum_{j \in \mathcal{N}_i} w_{ij} y(0) - \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_j^T \nabla y(0) - \frac{1}{2} \sum_{j \in \mathcal{N}_i} \tilde{x}_j^T (\mathcal{H}y)(0) \tilde{x}_j \\ &= -\frac{1}{2} \sum_{j \in \mathcal{N}_i} \tilde{x}_j^T (\mathcal{H}y)(0) \tilde{x}_j. \end{aligned}$$

If function $y(x)$ is a linear transform of the d -coordinates of x in the tangent space at x_i , then the Hessian matrix

$$(\mathcal{H}y)(0) := \left[\frac{\partial^2 y(x)}{\partial x(i) \partial x(j)} \right]_{x=0} = 0.$$

In this case $(I - W)y(0) = 0$ and y reaches a minimizer.

In other words, the kernel of Hessian operator \mathcal{H} has dimension $d+1$, consisting the constant function and d linearly independent coordinates. Inspired by such an observation, Donoho and Grimes [DG03b] proposed Hessian LLE (Eigenmap) in search of

$$\min_{y \perp \mathbf{1}} \int \|\mathcal{H}y\|^2, \quad \|y\| = 1.$$

The basic algorithmic idea is as follows.

1. G is incomplete, often k -nearest neighbour graph.
2. Local SVD on neighbourhood of x_i , for $x_{i,j} \in \mathcal{N}(x_i)$,

$$\tilde{X}^{(i)} = [x_{i_1} - \mu_i, \dots, x_{i_k} - \mu_i]^{p \times k} = \tilde{U}^{(i)} \tilde{\Sigma} (\tilde{V}^{(i)})^T,$$

where $\mu_i = \sum_{j=1}^k x_{i,j} = \frac{1}{k} X_i \mathbf{1}$. Here

- Left top singular vectors $\{\tilde{U}_1^{(i)}, \dots, \tilde{U}_d^{(i)}\}$ give an orthonormal basis of the approximate tangent space at x_i ,
- Right top singular vectors $[\tilde{V}_1^{(i)}, \dots, \tilde{V}_d^{(i)}]$ are representation coordinates in the tangent space of local sample points around x_i .

3. Null Hessian estimation: define

$$M = [1, \tilde{V}_1, \dots, \tilde{V}_d, \tilde{V}_1^2, \tilde{V}_1 \tilde{V}_2, \dots, \tilde{V}_{d-1} \tilde{V}_d, \tilde{V}_d^2] \in \mathbb{R}^{k \times (1+d+\binom{d+1}{2})}$$

where $\tilde{V}_i \tilde{V}_j = [\tilde{V}_{ik} \tilde{V}_{jk}]^T \in \mathbb{R}^k$ denotes the element-wise product (Hadamard product) between vector \tilde{V}_i and \tilde{V}_j .

Now we perform a Gram-Schmidt Orthogonalization procedure on M , get

$$\tilde{M} = [1, \hat{v}_1, \dots, \hat{v}_d, \hat{w}_1, \hat{w}_2, \dots, \hat{w}_{\binom{d+1}{2}}] \in \mathbb{R}^{k \times (1+d+\binom{d+1}{2})}$$

Define null Hessian by

$$[H^{(i)}]^T = [last \quad \binom{d+1}{2} \quad columns \quad of \quad \tilde{M}]_{k \times \binom{d+1}{2}},$$

as the first $d+1$ columns of \tilde{M} consists an orthonormal basis for the kernel of Hessian together with the constant vector.

Define a selection matrix $S^{(i)} \in \mathbb{R}^{n \times k}$ which selects those data in $\mathcal{N}(x_i)$, i.e.

$$[x_1, \dots, x_n]S^{(i)} = [x_{i_1}, \dots, x_{i_k}]$$

Then the kernel matrix is defined to be

$$K = \sum_{i=1}^n S^{(i)} H^{(i)T} H^{(i)} S^{(i)T} \in \mathbb{R}^{n \times n}$$

Find smallest $d+1$ eigenvectors of K and drop the smallest eigenvector, the remaining d eigenvectors will give rise to a d dimensional embedding of data points.

Algorithm 9: Hessian LLE Algorithm

Input: A weighted undirected graph $G = (V, E, d)$ such that

1 $V = \{x_i \in \mathbb{R}^p : i = 1, \dots, n\}$

2 $E = \{(i, j) : \text{if } j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors

Output: Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{d \times n}$ of data.

3 **Step 1:** Compute local PCA on neighborhood of x_i , for,

$$\tilde{X}^{(i)} = [x_{i_1} - \mu_i, \dots, x_{i_k} - \mu_i]^{p \times k} = \tilde{U}^{(i)} \tilde{\Sigma}(\tilde{V}^{(i)})^T, \quad x_{i_j} \in \mathcal{N}(x_i),$$

$$\text{where } \mu_i = \sum_{j=1}^k x_{i_j} = \frac{1}{k} X_i \mathbf{1};$$

4 **Step 2:** Null Hessian estimation: define

$$M = [1, \tilde{V}_1, \dots, \tilde{V}_d, \tilde{V}_1^2, \tilde{V}_1 \tilde{V}_2, \dots, \tilde{V}_{d-1} \tilde{V}_d, \tilde{V}_d^2] \in \mathbb{R}^{k \times (1+d+\binom{d+1}{2})}$$

where $\tilde{V}_i \tilde{V}_j = [\tilde{V}_{ik} \tilde{V}_{jk}]^T \in \mathbb{R}^k$ denotes the elementwise product (Hadamard product) between vector \tilde{V}_i and \tilde{V}_j . Now we perform a Gram-Schmidt Orthogonalization procedure on M , get

$$\tilde{M} = [1, \hat{v}_1, \dots, \hat{v}_d, \hat{w}_1, \hat{w}_2, \dots, \hat{w}_{\binom{d+1}{2}}] \in \mathbb{R}^{k \times (1+d+\binom{d+1}{2})}$$

Define

$$[H^{(i)}]^T = [last \quad \binom{d+1}{2} \quad columns \quad of \quad \tilde{M}]_{k \times \binom{d+1}{2}}.$$

Step 3: Define

$$K = \sum_{i=1}^n S^{(i)} H^{(i)T} H^{(i)} S^{(i)T} \in \mathbb{R}^{n \times n}, \quad [x_1, \dots, x_n]S^{(i)} = [x_{i_1}, \dots, x_{i_k}],$$

find smallest $d+1$ eigenvectors of K and drop the smallest eigenvector, and the remaining d eigenvectors will give rise to a d -embedding.

4.1. Convergence of Hessian LLE. There are two assumptions for the convergence of ISOMAP:

- Isometry: the geodesic distance between two points on manifolds equals to the Euclidean distances between intrinsic parameters.
- Convexity: the parameter space is a convex subset in \mathbb{R}^d .

Therefore, if the manifold contains a hole, ISOMAP will not faithfully recover the intrinsic coordinates. Hessian LLE above is provable to find local orthogonal coordinates for manifold reconstruction, even in nonconvex case. Figure [?] gives an example.

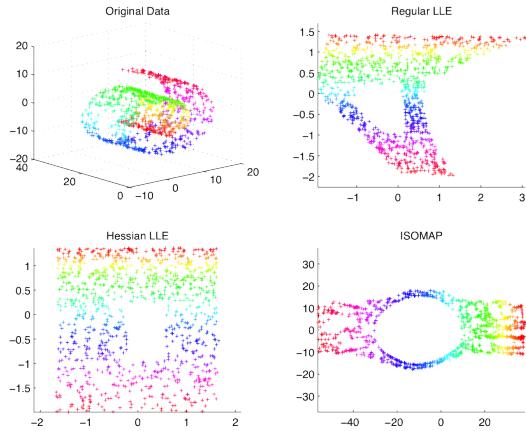


FIGURE 5. Comparisons of Hessian LLE on Swiss roll against ISOMAP and LLE. Hessian better recovers the intrinsic coordinates as the rectangular hole is the least distorted.

Donoho and Grimes [DG03b] relaxes the conditions above into the following ones.

- Local Isometry: in a small enough neighborhood of each point, geodesic distances between two points on manifolds are identical to Euclidean distances between parameter points.
- Connecteness: the parameter space is an open connected subset in \mathbb{R}^d .

Based on the relaxed conditions above, they prove the following result.

THEOREM 4.1. Supper $\mathcal{M} = \psi(\Theta)$ where Θ is an open connected subset of \mathbb{R}^d , and ψ is a locally isometric embedding of Θ into \mathbb{R}^n . Then the Hessian $\mathcal{H}(f)$ has a $d+1$ dimensional nullspace, consisting of the constant function and d -dimensional space of functions spanned by the original isometric coordinates.

Under this theorem, the original isometric coordinates can be recovered, up to a rigid motion, by identifying a suitable basis for the null space of $\mathcal{H}(f)$.

5. Local Tangent Space Alignment (LTSA)

A shortcoming of Hessian LLE lies in its bilinear form of local singular vectors (local PCA/MDS) to estimate the null Hessian. This is expensive when the intrinsic

dimensionality is high and also not stable when noise are presented. On the other hand, Zhenyue Zhang and Hongyuan Zha (2002) [ZZ02] suggest Local Tangent Space Alignment (LTSA) algorithm which just needs the linear form of local PCA which is more stable and cheaper than Hessian LLE.

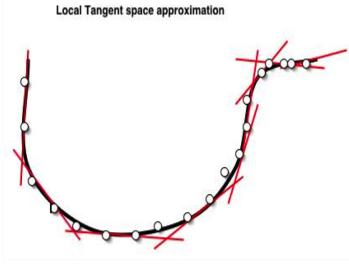


FIGURE 6. Local tangent space approximation.

The basic idea of LTSA is illustrated in Figure 6, where given a smooth curve (black), one can use discrete samples to find a good approximation of the tangent space of the original curve at each sample point. Finding such an approximation is in the spirit of *principal curve* or *principal manifold* proposed by Werner Stuetzle and Trevor Hastie [HS89]. Zhenyue Zhang and Hongyuan Zha (2002) [ZZ02] propose to use sampled data to find a good approximation of tangent space via local PCA, then the reconstruction data coordinates tries to preserve such approximate tangent space at each point to reach a global alignment.

Algorithm 10: LTSA Algorithm

Input: A weighted undirected graph $G = (V, E)$ such that

- 1 $V = \{x_i \in \mathbb{R}^p : i = 1, \dots, n\}$

- 2 $E = \{(i, j) : \text{if } j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors

Output: Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{k \times n}$ of data.

- 3 **Step 1** (local PCA): Compute local SVD on neighborhood of x_i , $x_{i_j} \in \mathcal{N}(x_i)$,

$$\tilde{X}^{(i)} = [x_{i_1} - \mu_i, \dots, x_{i_k} - \mu_i]^{p \times k} = \tilde{U}^{(i)} \tilde{\Sigma}(\tilde{V}^{(i)})^T,$$

where $\mu_i = \sum_{j=1}^k x_{i_j}$. Define

$$G_i = [1/\sqrt{k}, \tilde{V}_1^{(i)}, \dots, \tilde{V}_d^{(i)}]^{k \times (d+1)};$$

- 4 **Step 2** (tangent space alignment): Alignment (kernel) matrix

$$K^{n \times n} = \sum_{i=1}^n S_i W_i W_i^T S_i^T, \quad W_i^{k \times k} = I - G_i G_i^T,$$

where selection matrix $S_i^{n \times k} : [x_{i_1}, \dots, x_{i_k}] = [x_1, \dots, x_n] S_i^{n \times k}$;

- 5 **Step 3:** Find smallest $d + 1$ eigenvectors of K and drop the smallest eigenvector, the remaining d eigenvectors will give rise to a d -embedding.

For each x_i in \mathbb{R}^d with neighbor \mathcal{N}_i of size $|\mathcal{N}_i| = k_i - 1$, let $X^{(i)} = [x_{j_1}, x_{j_2}, \dots, x_{j_{k_i}}] \in \mathbb{R}^{p \times k_i}$ be the coordinate matrix. Consider the local SVD (PCA)

$$\tilde{X}^{(i)} = [x_{i_1} - \mu_i, \dots, x_{i_{k_i}} - \mu_i]^{p \times k_i} = X^{(i)} H = \tilde{U}^{(i)} \tilde{\Sigma}(\tilde{V}^{(i)})^T,$$

where $H = I - \frac{1}{k_i} \mathbf{1}_{k_i} \mathbf{1}_{k_i}^T$. Left singular vectors $\{\tilde{U}_1^{(i)}, \dots, \tilde{U}_d^{(i)}\}$ give an orthonormal basis of the approximate d -dimensional tangent space at x_i . Right singular vectors $(\tilde{V}_1^{(i)}, \dots, \tilde{V}_d^{(i)}) \cdot \tilde{\Sigma} \in \mathbb{R}^{k_i \times d}$ present the d -coordinates of k_i samples with respect to the tangent space basis.

Let $Y_i \in \mathbb{R}^{d \times k_i}$ be the embedding coordinates of the samples in \mathbb{R}^d and $L_i : \mathbb{R}^{p \times d}$ be an estimated basis of the tangent space at x_i in \mathbb{R}^p . Let $\Theta_i = \tilde{U}_d^{(i)} \tilde{\Sigma}_d (\tilde{V}_d^{(i)})^T \in \mathbb{R}^{p \times k_i}$ be the truncated SVD using top d components. LTSA looks for the minimizer of the following problem

$$(94) \quad \min_{Y, L} \sum_i \|E_i\|^2 = \sum_i \left\| Y_i \left(I - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) - L_i^T \Theta_i \right\|^2.$$

One can estimate $L_i^T = Y_i \left(I - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \Theta_i^\dagger$. Hence it reduces to

$$(95) \quad \min_Y \sum_i \|E_i\|^2 = \sum_i \left\| Y_i \left(I - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \left(I - \Theta_i^\dagger \Theta_i \right) \right\|^2$$

where $I - \Theta_i^\dagger \Theta_i$ is the projection to the normal space at x_i . This is equivalent to define

$$G_i = [1/\sqrt{k_i}, \tilde{V}_1^{(i)}, \dots, \tilde{V}_d^{(i)}]_{k_i \times (d+1)},$$

a weight matrix,

$$W_i^{k_i \times k_i} = I - G_i G_i^T,$$

and a positive semi-definite kernel matrix for alignment,

$$K^{n \times n} = \Phi = \sum_{i=1}^n S_i W_i W_i^T S_i^T$$

where the selection matrix $S_i^{n \times k_i} : [x_{i1}, \dots, x_{ik_i}] = [x_1, \dots, x_n] S_i$. Notice that constant vector is an eigenvector corresponding to the 0 eigenvalue. Hence similar to the LLE, one can choose bottom $d+1$ eigenvectors and drop the constant eigenvector, which gives embedding matrix $Y^{(n \times d)}$. An error analysis is given in [ZZ09], which shows that LTSA may recover the global coordinates asymptotically.

REMARK. We note that LTSA can be also applied to the situation that we are given local pairwise distances between samples. Since MDS and PCA are dual to each other, one can replace the local PCA in the algorithm by local MDS which leads to the same results as only right singular vectors $\tilde{V}^{(i)}$ are used there.

LTSA and Hessian LLE both may recover the linear coordinates, though the former is less expensive.

6. Laplacian LLE (Eigenmap)

Recall that in LLE, one chooses the weights w_{ij} to minimize the following energy

$$\min_{\sum_{j \in \mathcal{N}_i} w_{ij} = 1} \left\| \sum_{j \in \mathcal{N}_i} w_{ij} (x_j - x_i) \right\|^2.$$

In the ideal case, if the points $\tilde{x}_j = x_j - x_i$ are linearly dependent, then there is some w_{ij} , possibly not unique, such that $0 = \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_j$. Thus we have

$$0 = \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_j, \quad \text{and} \quad 1 = \sum_{j \in \mathcal{N}_i} w_{ij}.$$

For any smooth function $f(x)$, consider its Taylor expansion up to the second order

$$f(x) = f(0) + x^T \nabla f(0) + \frac{1}{2} x^T \mathcal{H}(0) x + o(\|x\|^2).$$

Therefore

$$\begin{aligned} (I - W)f(0) &:= f(0) - \sum_{j \in \mathcal{N}_i} w_{ij} f(\tilde{x}_i) \\ &\approx f(0) - \sum_{j \in \mathcal{N}_i} w_{ij} f(0) - \sum_{j \in \mathcal{N}_i} w_{ij} \tilde{x}_i^T \nabla f(0) - \frac{1}{2} \sum_{j \in \mathcal{N}_i} \tilde{x}_i^T \mathcal{H}(0) \tilde{x}_i \\ &= -\frac{1}{2} \sum_{j \in \mathcal{N}_i} \tilde{x}_i^T \mathcal{H}(0) \tilde{x}_i. \end{aligned}$$

When the $\{\tilde{x}_i\}$ in the last step becomes an orthonormal basis¹, the equation above gives

$$-\frac{1}{2} \sum_{j \in \mathcal{N}_i} \tilde{x}_i^T \mathcal{H}(0) \tilde{x}_i \approx \text{trace}(\mathcal{H}(0)) = \Delta f(0),$$

where the Laplacian operator $\Delta = \text{trace}(\mathcal{H}) = \sum_{i=1}^d \frac{\partial^2}{\partial \tilde{x}^2(i)}$. Such an observation leads to Laplacian LLE which looks for embedding functions

$$\min_{y \perp \mathbf{1}, \|y\|=1} \int \|\nabla y\|^2 = \int y^T \Delta y,$$

instead of in Hessian LLE,

$$\min_{y \perp \mathbf{1}, \|y\|=1} \int \|\mathcal{H}y\|^2.$$

The kernel of Laplacian consists of constant, linear functions, and bilinear functions of coordinates, of dimensionality $1 + d + \binom{d}{2}$. Therefore Laplacian LLE does not recover linear coordinates. However, Laplacian LLE converges to the spectrum of Laplacian-Beltrami operator, which enables us to choose w_{ij} as heat kernels. It has various connections with spectral graph theory and random walks on graphs, which further leads to *Diffusion Map* and relates to *topology* of data graph, namely the connectivity or the 0-th homology.

How to define Laplacian with discrete data? Graph Laplacians with heat kernels provide us an answer [BN01, BN03]. To see the idea, first consider a weighted oriented graph $G = (V, E, W)$ where $V = \{x_1, \dots, x_n\}$ is the vertex set, $E = \{(i, j) : i, j \in V\}$ is the set of oriented edges, and $W = [w_{ij} = w_{ji} \geq 0]$ is the weight matrix. Consider a particular weight matrix induced by heat kernels $W = (w_{ij}) \in \mathbb{R}^{n \times n}$ as

$$w_{ij} = \begin{cases} e^{-\frac{\|x_i - x_j\|^2}{t}} & j \in \mathcal{N}(i), \\ 0 & \text{otherwise.} \end{cases}$$

In particular, for $t \rightarrow \infty$, it gives binary weights. Let $D = \text{diag}(\sum_{j \in \mathcal{N}_i} w_{ij})$ be the diagonal matrix with weighted degree as diagonal elements. Define the *unnormalized graph Laplacian* by

$$L = D - W,$$

¹This is in general not true, which inspires Hessian LLE though.

and the *normalized graph Laplacian* by

$$\mathcal{L} = D^{-\frac{1}{2}}(D - W)D^{-\frac{1}{2}}.$$

To see the nature of such graph Laplacians, define the *graph gradient* map (also known as *co-boundary* map in algebraic topology) to be

$$\begin{aligned}\nabla : \mathbb{R}^V &\longrightarrow \mathbb{R}^E \\ f &\longmapsto (\nabla f)(i, j) = -(\nabla f)(j, i) = f(i) - f(j).\end{aligned}$$

Let $\widehat{D}_w = \text{diag}(w_{ij}) \in \mathbb{R}^{|E| \times |E|}$ be the diagonal matrix of edge weights. Then one can check that the unnormalized graph Laplacian satisfies $L = \nabla^T \widehat{D}_w \nabla$. In other words, it actually shows

$$f^T L f = (\nabla f)^T \widehat{D}_w (\nabla f) = \sum_{i \geq j} w_{ij} (f_i - f_j)^2.$$

This is the discrete analogue of the continuous Stokes Theorem on manifold,

$$\int \|\nabla_M f\|^2 = \int (\text{trace}(f^T \mathcal{H} f))^2,$$

where $\mathcal{H} = [\partial^2 / \partial_i \partial_j] \in \mathbb{R}^{d \times d}$ is the Hessian matrix.

For Laplacian Eigenmaps, there are two natural candidates, either the eigenvectors of unnormalized graph Laplacian L ,

$$\min_{y^T \mathbf{1}=0} \frac{y^T L y}{y^T y}$$

or generalized eigenvectors of L ,

$$\min_{y^T D \mathbf{1}=0} \frac{y^T L y}{y^T D y}.$$

A generalized eigenvector ϕ of L are also right eigenvectors of row Markov matrix $P = D^{-1}W$. To see this,

$$(D - W)\phi = \lambda D\phi \Leftrightarrow (I - D^{-1}W)v = \lambda v \Leftrightarrow D^{-1}W\phi = (1 - \lambda)\phi \Leftrightarrow P\phi = (1 - \lambda)\phi.$$

So eigenvectors are the same but only the eigenvalues are translated from λ to $1 - \lambda$.

In [BN03], it suggests *generalized eigenvectors* of L for Laplacian LLE, which scales the importance of vertex through its weighted degree and thus connects to random walk on graphs and diffusion map to be discussed later. Note that eigenvectors of normalized Laplacian \mathcal{L} are related to *generalized eigenvectors* of L up to a scaling matrix. This can be seen in the following reasoning.

$$\begin{aligned}\mathcal{L}v &= \lambda v \\ \Leftrightarrow D^{-\frac{1}{2}}(D - W)D^{-\frac{1}{2}}v &= \lambda v \\ \Leftrightarrow L\phi &= (D - W)\phi = \lambda D\phi, \quad \phi = D^{-\frac{1}{2}}v.\end{aligned}$$

In spectral graph theory, Fiedler theory actually tells us that the number of zero eigenvalues/generalized eigenvalues of L is the number of connected components of graph G (0-th Betti number); the corresponding eigenvectors can be used to partition the graph into components of small normalized cuts via Cheeger's inequality. On the other hand, lumpable Markov Chains on graphs will have piecewise constant Laplacian eigenmaps, which can be used for graph multiple normalized cut or partition.

Algorithm 11: Laplacian Eigenmap

-
- Input:** An adjacency graph $G = (V, E, d)$ such that
- 1 $V = \{x_i : i = 1, \dots, n\}$
 - 2 $E = \{(i, j) : \text{if } j \text{ is a neighbor of } i, \text{ i.e. } j \in \mathcal{N}_i\}$, e.g. k -nearest neighbors, ϵ -neighbors
 - 3 $d_{ij} = d(x_i, x_j)$, e.g. Euclidean distance for $x_i \sim x_j$ are in neighbor
- Output:** Euclidean d -dimensional coordinates $Y = [y_i] \in \mathbb{R}^{k \times n}$ of data.
- 4 **Step 1:** Choose weights
 - 5 (a) Heat kernel weights (parameter t):
- $$W_{ij} = \begin{cases} e^{-\frac{\|x_i - x_j\|^2}{t}}, & i \sim j, \\ 0, & \text{otherwise.} \end{cases}$$
- (b) Simple-minded ($t \rightarrow \infty$), $W_{ij} = 1$ if i and j are connected by an edge and $W_{ij} = 0$ otherwise.
- 6 **Step 2** (Eigenmap): Let $D = \text{diag}(\sum_j W_{ij})$ and $L = D - W$. Compute smallest $d + 1$ generalized eigenvectors
- $$Ly_l = \lambda_l Dy_l, \quad l = 0, 1, \dots, d,$$
- such that $0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_d$. Drop the zero eigenvalue λ_0 and constant eigenvector y_0 , and construct $Y_d = [y_1, \dots, y_d] \in \mathbb{R}^{n \times d}$.
-

To embed the data on to a d -dimensional Euclidean space, we can always choose bottom $d + 1$ eigenvectors, drop the smallest eigenvector (the constant vector associated with eigenvalue 0), and use the remaining d vectors to construct a d dimensional embedding of data.

6.1. Convergence of Laplacian Eigenmap. Some rigorous results about convergence of Laplacian eigenmaps are given in [BN08]. Assume that \mathcal{M} is a compact manifold with $\text{vol}(\mathcal{M}) = 1$. Let the Laplacian-Beltrami operator

$$\begin{aligned} \Delta_{\mathcal{M}} &: C(\mathcal{M}) \rightarrow L^2(\mathcal{M}) \\ f &\mapsto -\text{div}(\nabla f) \end{aligned}$$

Consider the following operator

$$\begin{aligned} \hat{L}_{t,n} &: C(\mathcal{M}) \rightarrow C(\mathcal{M}) \\ f &\mapsto \frac{1}{t(4\pi t)^{k/2}} \left(\sum_i e^{-\frac{\|y-x_i\|}{4t}} f(y) - \sum_i e^{\frac{\|y-x_i\|^2}{4t}} f(x_i) \right) \end{aligned}$$

where $(\hat{L}_{t,n}f)(y)$ is a function on \mathcal{M} , and

$$\begin{aligned} L_t &: L^2(\mathcal{M}) \rightarrow L^2(\mathcal{M}) \\ f &\mapsto \frac{1}{t(4\pi t)^{k/2}} \left(\int_{\mathcal{M}} e^{-\frac{\|y-x\|}{4t}} f(y) dx - \int_{\mathcal{M}} e^{\frac{\|y-x\|^2}{4t}} f(x) dx \right). \end{aligned}$$

Then [BN08] shows that when those operators have no repeated eigenvalues, the spectrum of $\hat{L}_{t,n}$ converges to L_t as $n \rightarrow \infty$ (variance), where the latter converges to that of $\Delta_{\mathcal{M}}$ with a suitable choice of $t \rightarrow \infty$ (bias). The following gives a summary.

THEOREM 6.1 (Belkin-Niyogi). Assume that all the eigenvalues in consideration are of multiplicity one. For small enough t , let $\hat{\lambda}_{n,i}^t$ be the i -th eigenvalue of $\hat{L}_{t,n}$

and $\hat{v}_{n,i}^t$ be the corresponding eigenfunction. Let λ_i and v_i be the corresponding eigenvalue and eigenfunction of $\Delta_{\mathcal{M}}$. Then there exists a sequence $t_n \rightarrow 0$ such that

$$\lim_{n \rightarrow \infty} \hat{\lambda}_{n,i}^{t_n} = \lambda_i$$

$$\lim_{n \rightarrow \infty} \|\hat{v}_{n,i}^{t_n} - v_i\| = 0$$

where the limits are taken in probability.

7. Diffusion Map

A detailed discussion on Diffusion Map will be after introducing random walks on graphs. In this section, we just make an introduction in comparison with Laplacian LLE.

Recall that for $x_i \in \mathbb{R}^d, i = 1, 2, \dots, n$, one can define a undirected weighted graph $G = (V, E, W)$ with $V = \{x_i : i = 1, \dots, n\}$, oriented edge set $E = \{(i, j)\}$, and symmetric weight $W = [w_{ij}]$ by heat kernel $w_{ij} = w_{ji} = \exp\left(-\frac{d(x_i, x_j)^2}{t}\right)$ for $i \sim j$ and $w_{ij} = 0$ otherwise. Assume that G is connected and thus has a finite diameter. Let $d_i = \sum_{j=1}^n W_{ij}$ and $D = \text{diag}(d_i)$.

A random walk on graph G can be defined through the following row Markov matrix,

$$P = D^{-1}W,$$

which is *primitive* (any two points can be connected by path of length no more than the diameter) and thus admits the following spectral decomposition

$$P = \Phi \Lambda \Psi^T,$$

where

- 1) $\Lambda = \text{diag}(\lambda_i)$ with $1 = \lambda_0 \geq \lambda_1 \geq \lambda_2 \dots \geq \lambda_{n-1} > -1$ for primitive Markov chain;
- 2) $\Phi = [\phi_0, \phi_1, \dots, \phi_{n-1}]$ are right eigenvectors of P , $P\Phi = \Phi\Lambda$;
- 3) $\Psi = [\psi_0, \psi_1, \dots, \psi_{n-1}]$ are left eigenvectors of P , $\Psi^T P = \Lambda \Psi^T$. Note that $\phi_0 = 1 \in \mathbb{R}^n$ and $\psi_0(i) = d_i / \sum_i d_i$. Thus ψ_0 is the same eigenvector as the stationary distribution $\pi(i) = d_i / \sum_i d_i$ ($\pi^T \mathbf{1} = \mathbf{1}$) up to a scaling factor;
- 4) Φ and Ψ are bi-orthogonal basis, *i.e.* $\phi_i^T \psi_j = \delta_{ij}$ or simply $\Phi^T \Psi = I$.

To see this, consider the normalized Laplacian

$$\mathcal{L} = D^{-1/2}(D - W)D^{-1/2},$$

which is symmetric and positive semi-definite, hence

$$S = D^{-1/2}WD^{-1/2} = I - \mathcal{L}$$

has n orthogonal eigenvectors $V = [v_1, v_2, \dots, v_n]$

$$S = V \Lambda V^T, \quad \Lambda = \text{diag}(\lambda_i), \quad V^T V = I,$$

where $1 = \lambda_0 \geq \lambda_1 \geq \lambda_2 \dots \geq \lambda_{n-1}$. Define $\Phi = D^{-1/2}V$ and $\Psi = D^{1/2}V$. One can obtain the spectral decomposition of P . Hence for any $\tau \geq 0$, $P^\tau = \Phi \Lambda^\tau \Psi^T$ defines a diffusion process on graph G , one can define a multiscale Euclidean embedding of data points.

Define diffusion map at scale t [CLL⁺05], by dropping the constant eigenvector ϕ_0 for connected graph G ,

$$\Phi_\tau(x_i) = [\lambda_1^\tau \phi_1(i), \dots, \lambda_{n-1}^\tau \phi_{n-1}(i)], \quad \tau \geq 0.$$

Clearly, Laplacian LLE corresponds to such a diffusion map at $\tau = 0$; as τ grows and small eigenvalues $|\lambda_i|^\tau < 1$ will drop to zero exponentially fast, which leads to a multiscale analysis on dimensionality reduction. For example, one can set a threshold $\delta > 0$, and only keep d_δ dimensions such that $|\lambda_i|^\tau \geq \delta$ for $1 \leq i \leq d_\delta$.

7.1. General Diffusion Maps and Convergence. In [CLL⁺05] a general class of diffusion maps are defined which involves a normalized weight matrix,

$$(96) \quad W_{ij}^{\alpha,t} = \frac{W_{ij}}{p_i^\alpha \cdot p_j^\alpha}, \quad p_i := \sum_k \exp\left(-\frac{d(x_i, x_k)^2}{t}\right)$$

where $\alpha = 0$ recovers the definition above. With this family, one can define $D_\alpha = \text{diag}(\sum_j W_{ij}^{\alpha,t})$ and the row Markov matrix

$$(97) \quad P_{\alpha,t,n} = D_\alpha^{-1} W^\alpha,$$

whose right eigenvectors Φ^α lead to a family of diffusion maps parameterized by α .

Such a definition suggests the following integral operators as diffusion operators. Assume that $q(x)$ is a density on \mathcal{M} .

- Let $k_t(x, y) = h(\|x - y\|^2/t)$ where h is a radial basis function, e.g. $h(z) = \exp(-z)$.
- Define

$$q_t(x) = \int_{\mathcal{M}} k_t(x, y) q(y) dy$$

and form the new kernel

$$k_t^{(\alpha)}(x, y) = \frac{k_t(x, y)}{q_t^\alpha(x) q_t^\alpha(y)}.$$

- Let

$$d_t^{(\alpha)}(x) = \int_{\mathcal{M}} k_t^{(\alpha)}(x, y) q(y) dy$$

and define the transition kernel of a Markov chain by

$$p_{t,\alpha}(x, y) = \frac{k_t^{(\alpha)}(x, y)}{d_t^{(\alpha)}(x)}.$$

Then the Markov chain can be defined as the operator

$$P_{t,\alpha} f(x) = \int_{\mathcal{M}} p_{t,\alpha}(x, y) f(y) q(y) dy.$$

- Define the infinitesimal generator of the Markov chain

$$L_{t,\alpha} = \frac{I - P_{t,\alpha}}{t}.$$

For this, Lafon et al.[CL06] shows the following pointwise convergence results.

THEOREM 7.1. Let $\mathcal{M} \in \mathbb{R}^p$ be a compact smooth submanifold, $q(x)$ be a probability density on \mathcal{M} , and $\Delta_{\mathcal{M}}$ be the Laplacian-Beltrami operator on \mathcal{M} .

$$(98) \quad \lim_{t \rightarrow 0} L_{t,\alpha} = \frac{\Delta_{\mathcal{M}}(f q^{1-\alpha})}{q^{1-\alpha}} - \frac{\Delta_{\mathcal{M}}(q^{1-\alpha})}{q^{1-\alpha}}.$$

This suggests that

- for $\alpha = 1$, it converges to the Laplacian-Beltrami operator $\lim_{t \rightarrow 0} L_{t,1} = \Delta_{\mathcal{M}}$;
- for $\alpha = 1/2$, it converges to a Schrödinger operator whose conjugation leads to a forward Fokker-Planck equation;
- for $\alpha = 0$, it is the normalized graph Laplacian.

A central question in Laplacian LLE and Diffusion maps is:

Why do we choose right eigenvectors ϕ_i of row Markov matrix for both Laplacian LLE and Diffusion map?

To answer this question, we will introduce Markov Chains on finite graphs to see various properties associated with their spectrum.

8. Stochastic Neighbor Embedding

Diffusion map preserves the diffusion distances

$$D^t(x_i, x_j) = \left(\sum_k \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2 \right)^{1/2} = \left(\sum_{k=1}^n \frac{(P(i,k) - P(j,k))^2}{d_k} \right)^{1/2},$$

like MDS. Such a diffusion distance is in fact a d_i^{-1} -weighted l_2 -distance between conditional probability representation of data points $P(i, \cdot)$ and $P(j, \cdot)$.

Instead of preserving diffusion distances, Stochastic Neighbor Embedding looks for embedding such that the estimated conditional probability $Q_Y(i, \cdot)$ is a faithful recover of P , by minimizing the Kullback-Leibler divergence between P and Q . P is similarly estimated using heat kernel as Diffusion maps. However, how to estimate Q_Y in a low-dimensional embedding space $Y = \mathbb{R}^d$? The original proposal of Stochastic Neighbor Embedding (SNE) uses the same heat kernel, which however suffers the crowding issue which push different classes of data points together. To overcome this issue, t-SNE exploits Student t -distribution or Cauchy distribution kernel which allows heavier tail than Gaussian distribution kernel, hence allowing more moderate distance points lying in the neighbor and attracting the same class members together.

to-be-finished...

9. Comparisons

According to the comparative studies by Todd Wittman, LTSA has the best overall performance in current manifold learning techniques. Try yourself his code, `mani.m`, and enjoy your new discoveries!

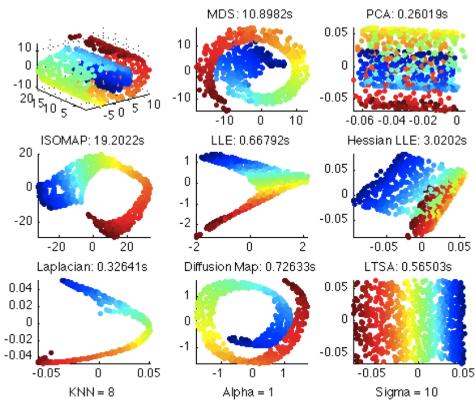


FIGURE 7. Comparisons of Manifold Learning Techniques on Swiss Roll

CHAPTER 6

Random Walk on Graphs

We have talked about Diffusion Map as a model of Random walk or Markov Chain on data graph. Among other methods of Manifold Learning, the distinct feature of Diffusion Map lies in that it combines both geometry and stochastic process. In the next few sections, we will talk about general theory of random walks or finite Markov chains on graphs which are related to data analysis. From this one can learn the origin of many ideas in diffusion maps.

Random Walk on Graphs.

- Perron-Frobenius Vector and Google’s PageRank: this is about Perron-Frobenius theory for nonnegative matrices, which leads to the characterization of nonnegative primary eigenvectors, such as stationary distributions of Markov chains; application examples include Google’s PageRank.
- Fiedler Vector, Cheeger’s Inequality, and Spectral Bipartition: this is about the second eigenvector in a Markov chain, mostly reduced from graph Laplacians (Fiedler theory, Cheeger’s Inequality), which is the basis for spectral partition.
- Lumpability/Metastability, piecewise constant right eigenvector, and Multiple spectral clustering (“MNCut” by Maila-Shi, 2001): this is about when to use multiple eigenvectors, whose relationship with lumpability or metastability of Markov chains, widely used in diffusion map, image segmentation, etc.
- Mean first passage time, commute time distance: the origins of diffusion distances.

Today we shall discuss the first part.

1. Introduction to Perron-Frobenius Theory and PageRank

Given $A_{n \times n}$, we define $A > 0$, *positive matrix*, iff $A_{ij} > 0 \forall i, j$, and $A \geq 0$, *nonnegative matrix*, iff $A_{ij} \geq 0 \forall i, j$.

Note that this definition is different from positive definite:

$$A \succ 0 \Leftrightarrow A \text{ is positive-definite} \Leftrightarrow x^T A x > 0 \quad \forall x \neq 0$$

$$A \succeq 0 \Leftrightarrow A \text{ is semi-positive-definite} \Leftrightarrow x^T A x \geq 0 \quad \forall x \neq 0$$

THEOREM 1.1 (Perron Theorem for Positive Matrix). Assume that $A > 0$, i.e.a positive matrix. Then

- 1) $\exists \lambda^* > 0, \nu^* > 0, \|\nu^*\|_2 = 1$, s.t. $A\nu^* = \lambda^*\nu^*$, ν^* is a right eigenvector
 $(\exists \lambda^* > 0, \omega > 0, \|\omega\|_2 = 1, \text{s.t. } (\omega^T)A = \lambda^*\omega^T)$, left eigenvector
- 2) \forall other eigenvalue λ of A , $|\lambda| < \lambda^*$
- 3) ν^* is unique up to rescaling or λ^* is simple

4) Collatz-Wielandt Formula

$$\lambda^* = \max_{x \geq 0, x \neq 0} \min_{x_i \neq 0} \frac{[Ax]_i}{x_i} = \min_{x > 0} \max \frac{[Ax]_i}{x_i}.$$

Such eigenvectors will be called Perron vectors. This result can be extended to nonnegative matrices.

THEOREM 1.2 (Nonnegative Matrix, Perron). Assume that $A \geq 0$, i.e. nonnegative. Then

- 1') $\exists \lambda^* > 0, \nu^* \geq 0, \|\nu^*\|_2 = 1$, s.t. $A\nu^* = \lambda^*\nu^*$ (similar to left eigenvector)
- 2') \forall other eigenvalue λ of A , $|\lambda| \leq \lambda^*$
- 3') ν^* is NOT unique

4) Collatz-Wielandt Formula

$$\lambda^* = \max_{x \geq 0, x \neq 0} \min_{x_i \neq 0} \frac{[Ax]_i}{x_i} = \min_{x > 0} \max \frac{[Ax]_i}{x_i}$$

Notice the changes in 1'), 2'), and 3'). Perron vectors are nonnegative rather than positive. In the nonnegative situation what we lose is the uniqueness in λ^* (2') and ν^* (3'). The next question is: can we add more conditions such that the loss can be remedied? Now recall the concept of irreducible and primitive matrices introduced before.

Irreducibility exactly describes the case that the induced graph from A is connected, i.e. every pair of nodes are connected by a path of arbitrary length. However primitivity strengthens this condition to k -connected, i.e. every pair of nodes are connected by a path of length k .

DEFINITION (Irreducible). The following definitions are equivalent:

- 1) For any $1 \leq i, j \leq n$, there is an integer $k \in \mathbb{Z}$, s.t. $A_{ij}^k > 0$; \Leftrightarrow
- 2) Graph $G = (V, E)$ ($V = \{1, \dots, n\}$ and $\{i, j\} \in E$ iff $A_{ij} > 0$) is (path-) connected, i.e. $\forall \{i, j\} \in E$, there is a path $(x_0, x_1, \dots, x_t) \in V^{n+1}$ where $i = x_0$ and $x_t = j$, connecting i and j .

DEFINITION (Primitive). The following characterizations hold:

- 1) There is an integer $k \in \mathbb{Z}$, such that $\forall i, j, A_{ij}^k > 0$; \Leftrightarrow
- 2) Any node pair $\{i, j\} \in E$ are connected with a path of length no more than k ; \Leftrightarrow
- 3) A has unique $\lambda^* = \max |\lambda|$; \Leftrightarrow
- 4) A is irreducible and $A_{ii} > 0$, for some i ,

Note that condition 4) is sufficient for primitivity but not necessary; all the first three conditions are necessary and sufficient for primitivity. Irreducible matrices have a simple primary eigenvalue λ^* and 1-dimensional primary (left and right) eigenspaces, with unique left and right eigenvectors. However, there might be other eigenvalues whose absolute values (module) equal to the primary eigenvalue, i.e., $\lambda^* e^{i\omega}$.

When A is a primitive matrix, A^k becomes a positive matrix for some k , then we can recover 1), 2) and 3) for positivity and uniqueness. This leads to the following Perron-Frobenius theorem.

THEOREM 1.3 (Nonnegative Matrix, Perron-Frobenius). Assume that $A \geq 0$ and A is primitive. Then

- 1) $\exists \lambda^* > 0, \nu^* > 0, \|\nu^*\|_2 = 1, s.t. A\nu^* = \lambda^*\nu^*$ (right eigenvector)
and $\exists \omega > 0, \|\omega\|_2 = 1, s.t. (\omega^T)A = \lambda^*\omega^T$ (left eigenvector)
- 2) \forall other eigenvalue λ of A , $|\lambda| < \lambda^*$
- 3) ν^* is unique
- 4) Collatz-Wielandt Formula

$$\lambda^* = \max_{x>0} \min_i \frac{[Ax]_i}{x_i} = \min_{x>0} \max_i \frac{[Ax]_i}{x_i}$$

Such eigenvectors and eigenvalue will be called as Perron-Frobenius or primary eigenvectors/eigenvalue.

EXAMPLE (Markov Chain). Given a graph $G = (V, E)$, consider a random walk on G with transition probability $P_{ij} = \text{Prob}(x_{t+1} = j | x_t = i) \geq 0$. Thus P is a row-stochastic or row-Markov matrix i.e. $P \cdot \vec{1} = \vec{1}$ where $\vec{1} \in \mathbb{R}^n$ is the vector with all elements being 1. From Perron theorem for nonnegative matrices, we know $\nu^* = \vec{1} > 0$ is a right Perron eigenvector of P
 $\lambda^* = 1$ is a Perron eigenvalue and all other eigenvalues $|\lambda| \leq 1 = \lambda^*$
 \exists left PF-eigenvector π such that $\pi^T P = \pi^T$ where $\pi \geq 0, \pi^T \pi = 1$; such π is called an invariant/equilibrium distribution
 P is irreducible (G is connected) $\Rightarrow \pi$ unique
 P is primitive (G connected by paths of length $\leq k$) $\Rightarrow |\lambda| = 1$ unique

$$\Leftrightarrow \lim_{t \rightarrow \infty} \pi_0^T P^k \rightarrow \pi^T \quad \forall \pi_0 \geq 0, \pi_0^T \pi_0 = 1$$

This means when we take powers of P , i.e. P^k , all rows of P^k will converge to the stationary distribution π^T . Such a convergence only holds when P is primitive. If P is not primitive, e.g. $P = [0, 1; 1, 0]$ (whose eigenvalues are 1 and -1), P^k always oscillates and never converges.

What's the rate of the convergence? Let

$$\gamma = \max\{|\lambda_2|, \dots, |\lambda_n|\}, \quad \lambda_1 = 1$$

and $\pi_t = (P^T)^t \pi_0$, roughly speaking we have

$$\|\pi_t - \pi\|_1 \sim O(e^{-\gamma t}).$$

This type of rates will be seen in various mixing time estimations.

A famous application of Markov chain in modern data analysis is Google's PageRank [BP98], although Google's current search engine only exploits that as one factor among many others. But you can still install Google Toolbar on your browser and inspect the PageRank scores of webpages. For more details about PageRank, readers may refer to Langville and Meyer's book [LM06].

EXAMPLE (PageRank). Consider a directed weighted graph $G = (V, E, W)$ whose weight matrix decodes the webpage link structure:

$$w_{ij} = \begin{cases} \#\{link : i \mapsto j\}, & (i, j) \in E \\ 0, & otherwise \end{cases}$$

Define an out-degree vector $d_i^o = \sum_{j=1}^n w_{ij}$, which measures the number of out-links from i . A diagonal matrix $D = \text{diag}(d_i)$ and a row Markov matrix $P_1 = D^{-1}W$, assumed for simplicity that all nodes have non-empty out-degree. This P_1 accounts

for a random walk according to the link structure of webpages. One would expect that stationary distributions of such random walks will disclose the importance of webpages: the more visits, the more important. However Perron-Frobenius above tells us that to obtain a unique stationary distribution, we need a primitive Markov matrix. For this purpose, Google's PageRank does the following trick.

Let $P_\alpha = \alpha P_1 + (1 - \alpha)E$, where $E = \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T$ is a random surfer model, *i.e.* one can jump to any other webpage uniformly. So in the model P_α , a browser will play a dice: he will jump according to link structure with probability α or randomly surf with probability $1 - \alpha$. With $1 > \alpha > 0$, the existence of random surfer model makes P a positive matrix, whence $\exists! \pi s.t. P_\alpha^T \pi = \pi$ (*means 'there exists a unique π '*). Google choose $\alpha = 0.85$ and in this case π gives PageRank scores.

Now you probably can figure out how to cheat PageRank. If there are many cross links between a small set of nodes (for example, Wikipedia), those nodes must appear to be high in PageRank. This phenomenon actually has been exploited by spam webpages, and even scholar citations. After learning the nature of PageRank, we should be aware of such mis-behaviors.

Finally we discussed a bit on Kleinberg's HITS algorithm [Kle99], which is based on singular value decomposition (SVD) of link matrix W . Above we have defined the out-degree d^o . Similarly we can define in-degree $d_k^i = \sum_j w_{jk}$. High out-degree webpages can be regarded as *hubs*, as they provide more links to others. On the other hand, high in-degree webpages are regarded as *authorities*, as they were cited by others intensively. Basically in/out-degrees can be used to rank webpages, which gives relative ranking as authorities/hubs. It turns out Kleinberg's HITS algorithm gives pretty similar results to in/out-degree ranking.

DEFINITION (HITS-authority). This use primary right singular vector of W as scores to give the ranking. To understand this, define $L_a = W^T W$. Primary right singular vector of W is just a primary eigenvector of nonnegative symmetric matrix L_a . Since $L_a(i, j) = \sum_k W_{ki} W_{kj}$, thus it counts the number of references which cites both i and j , *i.e.* $\sum_k \#\{i \leftarrow k \rightarrow j\}$. The higher value of $L_a(i, j)$ the more references received on the pair of nodes. Therefore Perron vector tend to rank the webpages according to authority.

DEFINITION (HITS-hub). This use primary left singular vector of W as scores to give the ranking. Define $L_h = W W^T$, whence primary left singular vector of W is just a primary eigenvector of nonnegative symmetric matrix L_h . Similarly $L_h(i, j) = \sum_k W_{ik} W_{jk}$, which counts the number of links from both i and j , hitting the same target, *i.e.* $\sum_k \#\{i \rightarrow k \leftarrow j\}$. Therefore the Perron vector L_h gives hub-ranking.

The last example is about Economic Growth model where the Debreu introduced nonnegative matrix into its study. Similar applications include population growth and exchange market, etc.

EXAMPLE (Economic Growth/Population/Exchange Market). Consider a market consisting n sectors (or families, currencies) where A_{ij} represents for each unit investment on sector j , how much the outcome in sector i . The nonnegative constraint $A_{ij} \geq 0$ requires that i and j are not *mutually inhibitor*, which means that investment in sector j does not decrease products in sector i . We study the dynamics $x_{t+1} = Ax_t$ and its long term behavior as $t \rightarrow \infty$ which describes the economic growth.

Moreover in exchange market, an additional requirement is put as $A_{ij} = 1/A_{ji}$, which is called *reciprocal matrix*. Such matrices are also used for preference aggregation in decision theory by Saaty.

From Perron-Frobenius theory we get: $\exists \lambda^* > 0 \quad \exists \nu^* \geq 0 \quad A\nu^* = \lambda^*\nu^*$ and $\exists \omega^* \geq 0 \quad A^T\omega^* = \lambda^*\omega^*$.

When A is primitive, ($A^k > 0$, i.e. investment in one sector will increase the product in another sector in no more than k industrial periods), we have for all other eigenvalues λ , $|\lambda| < \lambda^*$ and ω^*, ν^* are unique. In this case one can check that the long term economic growth is governed by

$$A^t \rightarrow (\lambda^*)^t \nu^* \omega^{*T}$$

where

- 1) for all i , $\frac{(x_t)_i}{(x_{t-1})_i} \rightarrow \lambda^*$
- 2) distribution of resources $\rightarrow \nu^*/\sum_i \nu_i^*$, so the distribution is actually not balanced
- 3) ω_i^* gives the relative value of investment on sector i in long term

1.1. Proof of Perron Theorem for Positive Matrices. A complete proof can be found in Meyer's book [Mey00], Chapter 8. Our proof below is based on optimization view, which is related to the Collatz-Wielandt Formula.

Assume that $A > 0$. Consider the following optimization problem:

$$\begin{aligned} & \max \delta \\ \text{s.t. } & Ax \geq \delta x \\ & x \geq 0 \\ & x \neq 0 \end{aligned}$$

Without loss of generality, assume that $1^T x = 1$. Let $y = Ax$ and consider the growth factor $\frac{y_i}{x_i}$, for $x_i \neq 0$. Our purpose above is to maximize the minimal growth factor δ ($y_i/x_i \geq \delta$).

Let λ^* be optimal value with $\nu^* \geq 0$, $1^T \nu^* = 1$, and $A\nu^* \geq \lambda^*\nu^*$. Our purpose is to show

- 1) $A\nu^* = \lambda^*\nu^*$
- 2) $\nu^* > 0$
- 3) ν^* and λ^* are unique.
- 4) For other eigenvalue λ ($\lambda z = Az$ when $z \neq 0$), $|\lambda| < \lambda^*$.

SKETCHY PROOF OF PERRON THEOREM. 1) If $A\nu^* \neq \lambda^*\nu^*$, then for some i , $[A\nu^*]_i > \lambda^*\nu_i^*$. Below we will find an increase of λ^* , which is thus not optimal. Define $\tilde{\nu} = \nu^* + \epsilon e_i$ with $\epsilon > 0$ and e_i denotes the vector which is one on the i^{th} component and zero otherwise.

For those $j \neq i$,

$$(A\tilde{\nu})_j = (A\nu^*)_j + \epsilon(Ae_i)_j = \lambda^*\nu_j^* + \epsilon A_{ji} > \lambda^*\nu_j^* = \lambda^*\tilde{\nu}_j$$

where the last inequality is due to $A > 0$.

For those $j = i$,

$$(A\tilde{\nu})_i = (A\nu^*)_i + \epsilon(Ae_i)_i > \lambda^*\nu_i^* + \epsilon A_{ii}.$$

Since $\lambda^* \tilde{\nu}_i = \lambda^* \nu_i^* + \epsilon \lambda^*$, we have

$$(A\tilde{\nu})_i - (\lambda^* \tilde{\nu})_i + \epsilon(A_{ii} - \lambda^*) = (A\nu^*)_i - (\lambda^* \nu_i^*) - \epsilon(\lambda^* - A_{ii}) > 0,$$

where the last inequality holds for small enough $\epsilon > 0$. That means, for some small $\epsilon > 0$, $(A\tilde{\nu}) > \lambda^* \tilde{\nu}$. Thus λ^* is not optimal, which leads to a contradiction.

2) Assume on the contrary, for some k , $\nu_k^* = 0$, then $(A\nu^*)_k = \lambda^* \nu_k^* = 0$. But $A > 0$, $\nu^* \geq 0$ and $\nu^* \neq 0$, so there $\exists i$, $\nu_i^* > 0$, which implies that $A\nu^* > 0$. That contradicts to the previous conclusion. So $\nu^* > 0$, which followed by $\lambda^* > 0$ (otherwise $A\nu^* > 0 = \lambda^* \nu^* = A\nu^*$).

3) We are going to show that for every $\nu \geq 0$, $A\nu = \mu\nu \Rightarrow \mu = \lambda^*$. Following the same reasoning above, A must have a left Perron vector $\omega^* > 0$, s.t. $A^T \omega^* = \lambda^* \omega^*$. Then $\lambda^*(\omega^{*T} \nu) = \omega^{*T} A\nu = \mu(\omega^{*T} \nu)$. Since $\omega^{*T} \nu > 0$ ($\omega^* > 0$, $\nu \geq 0$), there must be $\lambda^* = \mu$, i.e. λ^* is unique, and ν^* is unique.

4) For any other eigenvalue $Az = \lambda z$, $A|z| \geq |Az| = |\lambda||z|$, so $|\lambda| \leq \lambda^*$. Then we prove that $|\lambda| < \lambda^*$. Before proceeding, we need the following lemma.

LEMMA 1.4. $Az = \lambda z$, $|\lambda| = \lambda^*$, $z \neq 0 \Rightarrow A|z| = \lambda^*|z|$. $\lambda^* = \max_i |\lambda_i(A)|$

PROOF OF LEMMA. Since $|\lambda| = \lambda^*$,

$$A|z| = |A||z| \geq |Az| = |\lambda||z| = \lambda^*|z|$$

Assume that $\exists k$, $\frac{1}{\lambda^*} A|z_k| > |z_k|$. Denote $Y = \frac{1}{\lambda^*} A|z| - |z| \geq 0$, then $Y_k > 0$. Using that $A > 0$, $x \geq 0$, $x \neq 0 \Rightarrow Ax > 0$, we can get

$$\begin{aligned} &\Rightarrow \frac{1}{\lambda^*} AY > 0, \quad \frac{1}{\lambda^*} A|z| > 0 \\ &\Rightarrow \exists \epsilon > 0, \quad \frac{A}{\lambda^*} Y > \epsilon \frac{A}{\lambda^*} |z| \\ &\Rightarrow \bar{A}Y > \epsilon \bar{A}|z|, \quad \bar{A} = \frac{A}{\lambda^*} \\ &\Rightarrow \bar{A}^2|z| - \bar{A}|z| > \epsilon \bar{A}|z| \\ &\Rightarrow \frac{\bar{A}^2}{1+\epsilon}|z| > \bar{A}|z| \\ &\Rightarrow B = \frac{\bar{A}}{1+\epsilon}, \quad 0 = \lim_{m \rightarrow \infty} B^m \bar{A}|z| \geq \bar{A}|z| \\ &\Rightarrow \bar{A}|z| = 0 \quad \Rightarrow \quad |z| = 0 \quad \Rightarrow \quad Y = 0 \quad \Rightarrow \quad \bar{A}|z| = \lambda^*|z| \end{aligned}$$

□

Equipped with this lemma, assume that we have $Az = \lambda z$ ($z \neq 0$) with $|\lambda| = \lambda^*$, then

$$A|z| = \lambda^*|z| = |\lambda||z| = |Az| \quad \Rightarrow \quad \left| \sum_j \bar{a}_{ij} z_j \right| = \sum_j \bar{a}_{ij} |z_j|, \quad \bar{A} = \frac{A}{\lambda^*}$$

which implies that z_j has the same sign, i.e. $z_j \geq 0$ or $z_j \leq 0$ ($\forall j$). In both cases $|z|$ ($z \neq 0$) is a nonnegative eigenvector $A|z| = \lambda|z|$ which implies $\lambda = \lambda^*$ by 3). □

1.2. Perron-Frobenius theory for Nonnegative Tensors. Some researchers, e.g. Liqun Qi (Polytechnic University of Hong Kong), Lek-Heng Lim (U Chicago) and Kung-Ching Chang (PKU) et al. recently generalize Perron-Frobenius theory to nonnegative tensors, which may open a field toward *PageRank* for hypergraphs and array or tensor data. For example, $A(i, j, k)$ is a 3-tensor of dimension n , representing for each object $1 \leq i \leq n$, which object of j and k are closer to i .

A tensor of order- m and dimension- n means an array of n^m real numbers:

$$A = (a_{i_1, \dots, i_m}), \quad 1 \leq i_1, \dots, i_m \leq n$$

An n -vector $\nu = (\nu_1, \dots, \nu_n)^T$ is called an *eigenvector*, if

$$A\nu^{[m-1]} = \lambda\nu^{m-1}$$

for some $\lambda \in \mathbb{R}$, where

$$A\nu^{[m-1]} := \sum_{i_2, \dots, i_m=1}^n a_{k i_2 \dots i_m} \nu_{i_2} \cdots \nu_{i_m}, \quad \nu^{m-1} := (\nu_1^{m-1}, \dots, \nu_n^{m-1})^T.$$

Chang-Pearson-Zhang [2008] extends Perron-Frobenius theorem to show the existence of $\lambda^* > 0$ and $\nu^* > 0$ when $A > 0$ is irreducible.

$$\lambda^* = \max_{x>0} \min_i \frac{[Ax^{[m-1]}]_i}{x_i^{m-1}} = \min_{x>0} \max_i \frac{[Ax^{[m-1]}]_i}{x_i^{m-1}}.$$

2. Introduction to Fiedler Theory and Cheeger Inequality

In this class, we introduced the random walk on graphs. The last lecture shows Perron-Frobenius theory to the analysis of primary eigenvectors which is the stationary distribution. In this lecture we will study the second eigenvector. To analyze the properties of the graph, we construct two matrices: one is (unnormalized) graph Laplacian and the other is normalized graph Laplacian. In the first part, we introduce Fiedler Theory for the unnormalized graph Laplacian, which shows the second eigenvector can be used to bipartite the graph into two connected components. In the second part, we study the eigenvalues and eigenvectors of normalized Laplacian matrix to show its relations with random walks or Markov chains on graphs. In the third part, we will introduce the Cheeger Inequality for second eigenvector of normalized Laplacian, which leads to an approximate algorithm for Normalized graph cut (NCut) problem, an NP-hard problem itself.

2.1. Unnormalized Graph Laplacian and Fiedler Theory. Let $G = (V, E)$ be an undirected, unweighted simple¹ graph. Although the edges here are unweighted, the theory below still holds when weight is added. We can get a similar conclusion with the weighted adjacency matrix. However the extension to directed graphs will lead to different pictures.

We use $i \sim j$ to denote that node $i \in V$ is a neighbor of node $j \in V$.

DEFINITION (Adjacency Matrix).

$$A_{ij} = \begin{cases} 1 & i \sim j \\ 0 & \text{otherwise} \end{cases}.$$

¹Simple graph means for every pair of nodes there are at most one edge associated with it; and there is no self loop on each node.

REMARK. We can use the weight of edge $i \sim j$ to define A_{ij} if the graph is weighted. That indicates $A_{ij} \in \mathbb{R}^+$. We can also extend A_{ij} to \mathbb{R} which involves both positive and negative weights, like correlation graphs. But the theory below can not be applied to such weights being positive and negative.

The degree of node i is defined as follows.

$$d_i = \sum_{j=1}^n A_{ij}.$$

Define a diagonal matrix $D = \text{diag}(d_i)$. Now let's come to the definition of Laplacian Matrix L.

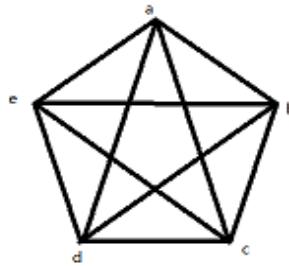
DEFINITION (Graph Laplacian).

$$L_{ij} = \begin{cases} d_i & i = j, \\ -1 & i \sim j \\ 0 & \text{otherwise} \end{cases}$$

This matrix is often called *unnormalized graph Laplacian* in literature, to distinguish it from the normalized graph Laplacian below. In fact, $L = D - A$.

EXAMPLE. $V = \{1, 2, 3, 4\}$, $E = \{\{1, 2\}, \{2, 3\}, \{3, 4\}\}$. This is a linear chain with four nodes.

$$L = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}.$$



EXAMPLE. A complete graph of n nodes, K_n . $V = \{1, 2, 3, \dots, n\}$, every two points are connected, as the figure above with $n = 5$.

$$L = \begin{pmatrix} n-1 & -1 & -1 & \dots & -1 \\ -1 & n-1 & -1 & \dots & -1 \\ -1 & \dots & -1 & n-1 & -1 \\ -1 & \dots & -1 & -1 & n-1 \end{pmatrix}.$$

From the definition, we can see that L is symmetric, so all its eigenvalues will be real and there is an orthonormal eigenvector system. Moreover L is positive semi-definite (p.s.d.). This is due to the fact that

$$\begin{aligned} v^T Lv &= \sum_i \sum_{j:j \sim i} v_i(v_i - v_j) = \sum_i \left(d_i v_i^2 - \sum_{j:j \sim i} v_i v_j \right) \\ &= \sum_{i \sim j} (v_i - v_j)^2 \geq 0, \quad \forall v \in \mathbb{R}^n. \end{aligned}$$

In fact, L admits the decomposition $L = BB^T$ where $B \in \mathbb{R}^{|V| \times |E|}$ is called *incidence matrix* (or *boundary map* in algebraic topology) here, for any $1 \leq j < k \leq n$,

$$B(i, \{j, k\}) = \begin{cases} 1, & i = j, \\ -1, & i = k, \\ 0, & \text{otherwise} \end{cases}$$

These two statements imply the eigenvalues of L can't be negative. That is to say $\lambda(L) \geq 0$.

THEOREM 2.1 (Fiedler theory). Let L has n eigenvectors

$$Lv_i = \lambda_i v_i, \quad v_i \neq 0, \quad i = 0, \dots, n-1$$

where $0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{n-1}$. For the second smallest eigenvector v_1 , define

$$N_- = \{i : v_1(i) < 0\},$$

$$N_+ = \{i : v_1(i) > 0\},$$

$$N_0 = V - N_- - N_+.$$

We have the following results.

- (1) $\#\{i, \lambda_i = 0\} = \#\{\text{connected components of } G\}$;
- (2) If G is connected, then both N_- and N_+ are connected. $N_- \cup N_0$ and $N_+ \cup N_0$ might be disconnected if $N_0 \neq \emptyset$.

This theorem tells us that the second smallest eigenvalue can be used to tell us if the graph is connected, *i.e.* G is connected iff $\lambda_1 \neq 0$, *i.e.*

$$\lambda_1 = 0 \Leftrightarrow \text{there are at least two connected components.}$$

$$\lambda_1 > 0 \Leftrightarrow \text{the graph is connected.}$$

Moreover, the second smallest eigenvector can be used to bipartite the graph into two connected components by taking N_- and N_+ when N_0 is empty. For this reason, we often call the second smallest eigenvalue λ_1 as the *algebraic connectivity*. More materials can be found in Jim Demmel's Lecture notes on Fiedler Theory at UC Berkeley: why we use unnormalized Laplacian eigenvectors for spectral partition (<http://www.cs.berkeley.edu/~demmel/cs267/lecture20/lecture20.html>).

We can calculate eigenvalues by using Rayleigh Quotient. This gives a sketch proof of the first part of the theory.

PROOF OF PART I. Let (λ, v) be a pair of eigenvalue-eigenvector, *i.e.* $Lv = \lambda v$. Since $L1 = 0$, so the constant vector $1 \in \mathbb{R}^n$ is always the eigenvector associated with $\lambda_0 = 0$. In general,

$$\lambda = \frac{v^T Lv}{v^T v} = \frac{\sum_{i \sim j} (v_i - v_j)^2}{\sum_i v_i^2}.$$

Note that

$$0 = \lambda_1 \Leftrightarrow v_i = v_j \text{ (} j \text{ is path connected with } i \text{).}$$

Therefore v is a piecewise constant function on connected components of G . If G has k components, then there are k independent piecewise constant vectors in the span of characteristic functions on those components, which can be used as eigenvectors of L . In this way, we proved the first part of the theory. \square

2.2. Normalized graph Laplacian and Cheeger's Inequality.

DEFINITION (Normalized Graph Laplacian).

$$\mathcal{L}_{ij} = \begin{cases} 1 & i = j, \\ -\frac{1}{\sqrt{d_i d_j}} & i \sim j, \\ 0 & \text{otherwise.} \end{cases}$$

In fact $\mathcal{L} = D^{-1/2}(D - A)D^{-1/2} = D^{-1/2}LD^{-1/2} = I - D^{-1/2}(D - A)D^{-1/2}$. From this one can see the relations between eigenvectors of normalized \mathcal{L} and unnormalized L . For eigenvectors $\mathcal{L}v = \lambda v$, we have

$$(I - D^{-1/2}LD^{-1/2})v = \lambda v \Leftrightarrow Lu = \lambda Du, \quad u = D^{-1/2}v,$$

whence eigenvectors of \mathcal{L} , v after rescaling by $D^{-1/2}v$, become generalized eigenvectors of L .

We can also use the Rayleigh Quotient to calculate the eigenvalues of \mathcal{L} .

$$\begin{aligned} \frac{v^T \mathcal{L}v}{v^T v} &= \frac{v^T D^{-\frac{1}{2}}(D - A)D^{-\frac{1}{2}}v}{v^T v} \\ &= \frac{u^T Lu}{u^T Du} \\ &= \frac{\sum_{i \sim j} (u_i - u_j)^2}{\sum_j u_j^2 d_j}. \end{aligned}$$

Similarly we get the relations between eigenvalue and the connected components of the graph.

$$\#\{\lambda_i(\mathcal{L}) = 0\} = \#\{\text{connected components of } G\}.$$

Next we show that eigenvectors of \mathcal{L} are related to random walks on graphs. This will show you why we choose this matrix to analysis the graph.

We can construct a random walk on G whose transition matrix is defined by

$$P_{ij} \sim \frac{A_{ij}}{\sum_j A_{ij}} = \frac{1}{d_i}.$$

By easy calculation, we see the result below.

$$P = D^{-1}A = D^{-1/2}(I - \mathcal{L})D^{1/2}.$$

Hence P is similar to $I - \mathcal{L}$. So their eigenvalues satisfy $\lambda_i(P) = 1 - \lambda_i(\mathcal{L})$. Consider the right eigenvector ϕ and left eigenvector ψ of P .

$$u^T P = \lambda u,$$

$$Pv = \lambda v.$$

Due to the similarity between P and \mathcal{L} ,

$$u^T P = \lambda u^T \Leftrightarrow u^T D^{-1/2} (I - \mathcal{L}) D^{1/2} = \lambda u^T.$$

Let $\bar{u} = D^{-1/2} u$, we will get:

$$\begin{aligned} \bar{u}^T (I - \mathcal{L}) &= \lambda \bar{u}^T \\ \Leftrightarrow \mathcal{L} \bar{u} &= (1 - \lambda) \bar{u}. \end{aligned}$$

You can see \bar{u} is the eigenvector of \mathcal{L} , and we can get left eigenvectors of P from \bar{u} by multiply it with $D^{1/2}$ on the left side. Similarly for the right eigenvectors $v = D^{-1/2} \bar{u}$.

If we choose $u_0 = \pi_i \sim \frac{d_i}{\sum d_i}$, then:

$$\begin{aligned} \bar{u}_0(i) &\sim \sqrt{d_i}, \\ \bar{u}_k^T \bar{u}_l &= \delta_{kl}, \\ u_k^T D v_l &= \delta_{kl}, \\ \pi_i P_{ij} &= \pi_j P_{ji} \sim A_{ij} = A_{ji}, \end{aligned}$$

where the last identity says the Markov chain is time-reversible.

All the conclusions above show that the normalized graph Laplacian \mathcal{L} keeps some connectivity measure of unnormalized graph Laplacian L . Furthermore, \mathcal{L} is more related with random walks on graph, through which eigenvectors of P are easy to check and calculate. That's why we choose this matrix to analysis the graph.

Now we are ready to introduce the Cheeger's inequality with normalized graph Laplacian.

Let G be a graph, $G = (V, E)$ and S is a subset of V whose complement is $\bar{S} = V - S$. We define $Vol(S)$, $CUT(S)$ and $NCUT(S)$ as below.

$$Vol(S) = \sum_{i \in S} d_i.$$

$$CUT(S) = \sum_{i \in S, j \in \bar{S}} A_{ij}.$$

$$NCUT(S) = \frac{CUT(S)}{\min(Vol(S), Vol(\bar{S}))}.$$

$NCUT(S)$ is called normalized-cut. We define the Cheeger constant

$$h_G = \min_S NCUT(S).$$

Finding minimal normalized graph cut is NP-hard. It is often defined that

$$\text{Cheeger ratio (expander): } h_S := \frac{CUT(S)}{Vol(S)}$$

and

$$\text{Cheeger constant: } h_G := \min_S \max \{h_S, h_{\bar{S}}\}.$$

Cheeger Inequality says the second smallest eigenvalue provides both upper and lower bounds on the minimal normalized graph cut. Its proof gives us a constructive polynomial algorithm to achieve such bounds.

THEOREM 2.2 (Cheeger Inequality). For every undirected graph G ,

$$\frac{h_G^2}{2} \leq \lambda_1(\mathcal{L}) \leq 2h_G.$$

PROOF. (1) Upper bound:

Assume the following function f realizes the optimal normalized graph cut,

$$f(i) = \begin{cases} \frac{1}{Vol(S)} & i \in S, \\ -\frac{1}{Vol(\bar{S})} & i \in \bar{S}, \end{cases}$$

By using the Rayleigh Quotient, we get

$$\begin{aligned} \lambda_1 &= \inf_{g \perp D^{1/2} e} \frac{g^T \mathcal{L} g}{g^T g} \\ &\leq \frac{\sum_{i \sim j} (f_i - f_j)^2}{\sum f_i^2 d_i} \\ &= \frac{(\frac{1}{Vol(S)} + \frac{1}{Vol(\bar{S})})^2 CUT(S)}{Vol(S) \frac{1}{Vol(S)^2} + Vol(\bar{S}) \frac{1}{Vol(\bar{S})^2}} \\ &= (\frac{1}{Vol(S)} + \frac{1}{Vol(\bar{S})}) CUT(S) \\ &\leq \frac{2CUT(S)}{\min(Vol(S), Vol(\bar{S}))} =: 2h_G. \end{aligned}$$

which gives the upper bound.

(2) Lower bound: the proof of lower bound actually gives a constructive algorithm to compute an approximate optimal cut as follows.

Let v be the second eigenvector, i.e. $\mathcal{L}v = \lambda_1 v$, and $f = D^{-1/2}v$. Then we reorder node set V such that $f_1 \leq f_2 \leq \dots \leq f_n$. Denote $V_- = \{i; v_i < 0\}$, $V_+ = \{i; v_i \geq v_r\}$. Without Loss of generality, we can assume

$$\sum_{i \in V_-} d_v \geq \sum_{i \in V_+} d_v$$

Define new functions f^+ to be the magnitudes of f on V_+ .

$$f_i^+ = \begin{cases} f_i & i \in V_+, \\ 0 & otherwise, \end{cases}$$

Now consider a series of particular subsets of V ,

$$S_i = \{v_1, v_2, \dots, v_i\},$$

and define

$$\widehat{Vol}(S) = \min(Vol(S), Vol(\bar{S})).$$

$$\alpha_G = \min_i NCUT(S_i).$$

Clearly finding the optimal value α just requires comparison over $n - 1$ NCUT values.

Below we shall show that

$$\frac{h_G^2}{2} \leq \frac{\alpha_G^2}{2} \leq \lambda_1.$$

First, we have $Lf = \lambda_1 Df$, so we must have

$$(99) \quad \sum_{j:j \sim i} f_i(f_i - f_j) = \lambda_1 d_i f_i^2.$$

From this we will get the following results,

$$\begin{aligned} \lambda_1 &= \frac{\sum_{i \in V_+} f_i \sum_{j:j \sim i} (f_i - f_j)}{\sum_{i \in V_+} d_i f_i^2}, \\ &= \frac{\sum_{i \sim j} i, j \in V_+ (f_i - f_j)^2 + \sum_{i \in V_+} f_i \sum_{j \sim i} j \in V_- (f_i - f_j)}{\sum_{i \in V_+} d_i f_i^2}, (f_i - f_j)^2 = f_i(f_i - f_j) + f_j(f_j - f_i) \\ &> \frac{\sum_{i \sim j} i, j \in V_+ (f_i - f_j)^2 + \sum_{i \in V_+} f_i \sum_{j \sim i} j \in V_- (f_i)}{\sum_{i \in V_+} d_i f_i^2}, \\ &= \frac{\sum_{i \sim j} (f_i^+ - f_j^+)^2}{\sum_{i \in V} d_i f_i^{+2}}, \\ &= \frac{(\sum_{i \sim j} (f_i^+ - f_j^+)^2)(\sum_{i \sim j} (f_i^+ + f_j^+)^2)}{(\sum_{i \in V} f_i^{+2} d_i)(\sum_{i \sim j} (f_i^+ + f_j^+)^2)} \\ &\geq \frac{(\sum_{i \sim j} f_i^{+2} - f_j^{+2})^2}{(\sum_{i \in V} f_i^{+2} d_i)(\sum_{i \sim j} (f_i^+ + f_j^+)^2)}, \text{ Cauchy-Schwartz Inequality} \\ &\geq \frac{(\sum_{i \sim j} f_i^{+2} - f_j^{+2})^2}{2(\sum_{i \in V} f_i^{+2} d_i)^2}, \end{aligned}$$

where the second last step is due to the Cauchy-Schwartz inequality $|\langle x, y \rangle|^2 \leq \langle x, x \rangle \cdot \langle y, y \rangle$, and the last step is due to $\sum_{i \sim j \in V} (f_i^+ + f_j^+)^2 = \sum_{i \sim j \in V} (f_i^{+2} + f_j^{+2} + 2f_i^+ f_j^+) \leq 2 \sum_{i \sim j \in V} (f_i^{+2} + f_j^{+2}) \leq 2 \sum_{i \in V} f_i^{+2} d_i$. Continued from the last inequality,

$$\begin{aligned} \lambda_1 &\geq \frac{(\sum_{i \sim j} f_i^{+2} - f_j^{+2})^2}{2(\sum_{i \in V} f_i^{+2} d_i)^2}, \\ &\geq \frac{(\sum_{i \in V} (f_i^{+2} - f_{i-1}^{+2}) \text{CUT}(S_{i-1}))^2}{2(\sum_{i \in V} f_i^{+2} d_i)^2}, \text{ since } f_1 \leq f_2 \leq \dots \leq f_n \\ &\geq \frac{(\sum_{i \in V} (f_i^{+2} - f_{i-1}^{+2}) \alpha_G \widetilde{\text{Vol}}(S_{i-1}))^2}{2(\sum_{i \in V} f_i^{+2} d_i)^2} \\ &= \frac{\alpha_G^2}{2} \cdot \frac{(\sum_{i \in V} f_i^{+2} (\widetilde{\text{Vol}}(S_{i-1}) - \widetilde{\text{Vol}}(S_i)))^2}{(\sum_{i \in V} f_i^{+2} d_i)^2}, \\ &= \frac{\alpha_G^2}{2} \frac{(\sum_{i \in V} f_i^{+2} d_i)^2}{(\sum_{i \in V} f_i^{+2} d_i)^2} = \frac{\alpha_G^2}{2}. \end{aligned}$$

where the last inequality is due to the assumption $\text{Vol}(V_-) \geq \text{Vol}(V_+)$, whence $\widetilde{\text{Vol}}(S_i) = \text{Vol}(\bar{S}_i)$ for $i \in V_+$.

This completes the proof. \square

Fan Chung gives a short proof of the lower bound in Simons Institute workshop, 2014.

SHORT PROOF. The proof is based on the fact that

$$h_G = \inf_{f \neq 0} \sup_{c \in \mathbb{R}} \frac{\sum_{x \sim y} |f(x) - f(y)|}{\sum_x |f(x) - c| d_x}$$

where the supreme over c is reached at $c^* = \text{median}(f(x) : x \in V)$.

$$\begin{aligned} \lambda_1 &= R(f) = \sup_c \frac{\sum_{x \sim y} (f(x) - f(y))^2}{\sum_x (f(x) - c)^2 d_x}, \\ &\geq \frac{\sum_{x \sim y} (g(x) - g(y))^2}{\sum_x g(x)^2 d_x}, \quad g(x) = f(x) - c \\ &= \frac{(\sum_{x \sim y} (g(x) - g(y))^2)(\sum_{x \sim y} (g(x) + g(y))^2)}{(\sum_{x \in V} g^2(x) d_x)((\sum_{x \sim y} (g(x) + g(y))^2))} \\ &\geq \frac{(\sum_{x \sim y} |g^2(x) - g^2(y)|)^2}{(\sum_{x \in V} g^2(x) d_x)((\sum_{x \sim y} (g(x) + g(y))^2))}, \quad \text{Cauchy-Schwartz Inequality} \\ &\geq \frac{(\sum_{x \sim y} |g^2(x) - g^2(y)|)^2}{2(\sum_{x \in V} g^2(x) d_x)^2}, \quad (g(x) + g(y))^2 \leq 2(g^2(x) + g^2(y)) \\ &\geq \frac{h_G^2}{2}. \end{aligned}$$

\square

3. *Laplacians and the Cheeger inequality for directed graphs

The following section is mainly contained in [Chu05], which described the following results:

- (1) Define Laplacians on directed graphs.
- (2) Define Cheeger constants on directed graphs.
- (3) Give an example of the singularity of Cheeger constant on directed graph.
- (4) Use the eigenvalue of Lapacian and the Cheeger constant to estimate the convergence rate of random walk on a directed graph.

Another good reference is [LZ10].

3.1. Definition of Laplacians on directed graphs. On a finite and strong connected directed graph $G = (V, E)$ (A directed graph is strong connected if there is a path between any pair of vertices), a weight is a function

$$w : E \rightarrow \mathbb{R}_{\geq 0}$$

The in-degree and out-degree of a vertex are defined as

$$\begin{aligned} d^{in} : V &\rightarrow \mathbb{R}_{\geq 0} \\ d_i^{in} &= \sum_{j \in V} w_{ji} \\ d^{out} : V &\rightarrow \mathbb{R}_{\geq 0} \\ d_i^{out} &= \sum_{j \in V} w_{ij} \end{aligned}$$

Note that d_i^{in} may be different from d_i^{out} .

A random walk on the weighted G is a Markov chain with transition probability

$$P_{ij} = \frac{w_{ij}}{d_i^{out}}.$$

Since G is strong connected, P is irreducible, and consequently there is a unique stationary distribution ϕ . (And the distribution of the Markov chain will converge to it if and only if P is aperiodic.)

EXAMPLE (undirected graph).

$$\phi(x) = \frac{d_x}{\sum_y d_y}.$$

EXAMPLE (Eulerian graph). If $d_x^{in} = d_x^{out}$ for every vertex x , then $\phi(x) = \frac{d_x^{out}}{\sum_y d_y^{out}}$.

This is because d_x^{out} is an unchanged measure with

$$\sum_x d_x^{out} P_{xy} = \sum_x w_{xy} = d_y^{in} = d_y^{out}.$$

EXAMPLE (exponentially small stationary dist.). G is a directed graph with $n + 1$ vertices formed by the union of a directed circle $v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n$ and edges $v_i \rightarrow v_0$ for $i = 1, 2, \dots, n$. The weight on any edge is 1. Checking from v_n to v_0 with the prerequisite of stationary distribution that the inward probability flow equals to the outward probability flow, we can see that

$$\phi(v_0) = 2^n \phi(v_n), \text{ i.e. } \phi(v_n) = 2^{-n} \phi(v_0).$$

This exponentially small stationary distribution cannot occur in undirected graph cases for then

$$\phi(i) = \frac{d_i}{\sum_j d_j} \geq \frac{1}{n(n-1)}.$$

However, the stationary dist. can be no smaller than exponential, because we have

THEOREM 3.1. If G is a strong connected directed graph with $w \equiv 1$, and $d_x^{out} \leq k, \forall x$, then $\max\{\phi(x) : x \in V\} \leq k^D \min\{\phi(y) : y \in V\}$, where D is the diameter of G .

It can be easily proved using induction on the path connecting x and y .

Now we give a definition on those balanced weights.

DEFINITION (circulation).

$$F : E \rightarrow \mathbb{R}_{\geq 0}$$

If F satisfies

$$\sum_{u,u \rightarrow v} F(u,v) = \sum_{w,v \rightarrow w} F(v,w), \forall v,$$

then F is called a circulation.

NOTE. A circulation is a flow with no source or sink.

EXAMPLE. For a directed graph, $F_\phi(u, v) = \phi(u)P(u, v)$ is a circulation, for

$$\sum_{u, u \rightarrow v} F_\phi(u, v) = \phi(v) = \sum_{w, v \rightarrow w} F_\phi(v, w).$$

DEFINITION (Rayleigh quotient). For a directed graph G with transition probability matrix P and stationary distribution ϕ , the Rayleigh quotient for any $f : V \rightarrow \mathbb{C}$ is defined as

$$R(f) = \frac{\sum_{u \rightarrow v} |f(u) - f(v)|^2 \phi(u)P(u, v)}{\sum_v |f(v)|^2 \phi(v)}.$$

NOTE. Compare with the undirected graph condition where

$$R(f) = \frac{\sum_{u \sim v} |f(u) - f(v)|^2 w_{uv}}{\sum_v |f(v)|^2 d(v)}.$$

If we look on every undirected edge (u, v) as two directed edges $u \rightarrow v, v \rightarrow u$, then we get a Eulerian directed graph. So $\phi(u) \sim d_u^{out}$ and $d_u^{out}P(u, v) = w_{uv}$, as a result $R(f)(directed) = 2R(f)(undirected)$. The factor 2 is the result of looking on every edge as two edges.

The next step is to extend the definition of Laplacian to directed graphs. First we give a review on Laplacian on undirected graphs. On an undirected graph, adjacent matrix is

$$\begin{aligned} A_{ij} &= \begin{cases} 1, & i \sim j; \\ 0, & i \not\sim j. \end{cases} \\ D &= \text{diag}(d(i)), \\ \mathcal{L} &= D^{-1/2}(D - A)D^{-1/2}. \end{aligned}$$

On a directed graph, however, there are two degrees on a vertex which are generally inequivalent. Notice that on an undirected graph, stationary distribution $\phi(i) \sim d(i)$, so $D = c\Phi$, where c is a constant and $\Phi = \text{diag}(\phi(i))$.

$$\begin{aligned} \mathcal{L} &= I - D^{-1/2}AD^{-1/2} \\ &= I - D^{1/2}PD^{-1/2} \\ &= I - c^{1/2}\Phi^{1/2}Pc^{-1/2}\Phi^{-1/2} \\ &= I - \Phi^{1/2}P\Phi^{-1/2} \end{aligned}$$

Extending and symmetrizing it, we define Laplacian on a directed graph

DEFINITION (Laplacian).

$$\mathcal{L} = I - \frac{1}{2}(\Phi^{1/2}P\Phi^{-1/2} + \Phi^{-1/2}P^*\Phi^{1/2}).$$

Suppose the eigenvalues of \mathcal{L} are $0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{n-1}$. Like the undirected case, we can calculate λ_1 with the Rayleigh quotient.

THEOREM 3.2.

$$\lambda_1 = \inf_{\sum f(x)\phi(x)=0} \frac{R(f)}{2}.$$

Before proving that, we need

LEMMA 3.3.

$$R(f) = 2 \frac{g\mathcal{L}g^*}{\|g\|^2}, \text{ where } g = f\Phi^{1/2}.$$

PROOF.

$$\begin{aligned}
R(f) &= \frac{\sum_{u \rightarrow v} |f(u) - f(v)|^2 \phi(u) P(u, v)}{\sum_v |f(v)|^2 \phi(v)} \\
&= \frac{\sum_{u \rightarrow v} |f(u)|^2 \phi(u) P(u, v) + \sum_v |f(v)|^2 \phi(v) - \sum_{u \rightarrow v} (\overline{f(u)} f(v) + f(u) \overline{f(v)}) \phi(u) P(u, v)}{f \Phi f^*} \\
&= \frac{\sum_u |f(u)|^2 \phi(u) + \sum_v |f(v)|^2 \phi(v) - (f^* \Phi P f + f \Phi P f^*)}{f \Phi f^*} \\
&= 2 - \frac{f(P^* \Phi + \Phi P) f^*}{f \Phi f^*} \\
&= 2 - \frac{(g \Phi^{-1/2})(P^* \Phi + \Phi P)(\Phi^{-1/2} g^*)}{(g \Phi^{-1/2}) \Phi(\Phi^{-1/2} g^*)} \\
&= 2 - \frac{g(\Phi^{-1/2} P^* \Phi^{1/2} + \Phi^{1/2} P \Phi^{-1/2}) g^*}{g g^*} \\
&= 2 \cdot \frac{g \mathcal{L} g^*}{\|g\|^2}
\end{aligned}$$

□

PROOF OF THEOREM 3.2. With Lemma 3.3 and $\mathcal{L}(\phi(x)^{1/2})_{n \times 1} = 0$, we have

$$\begin{aligned}
\lambda_1 &= \inf_{\sum g(x) \phi(x)^{1/2} = 0} \frac{R(f)}{2} \\
&= \inf_{\sum f(x) \phi(x) = 0} \frac{R(f)}{2}.
\end{aligned}$$

□

NOTE.

$$\begin{aligned}
\lambda_1 &= \inf_{f, \sum f(x) \phi(x) = 0} \frac{R(f)}{2} \\
&= \inf_{f, \sum f(x) \phi(x) = 0} \frac{\sum_{u \rightarrow v} |f(u) - f(v)|^2 \phi(u) P(u, v)}{2 \sum_v |f(v)|^2 \phi(v)} \\
&= \inf_{f, \sum f(x) \phi(x) = 0} \sup_c \frac{\sum_{u \rightarrow v} |f(u) - f(v)|^2 \phi(u) P(u, v)}{2 \sum_v |f(v) - c|^2 \phi(v)}
\end{aligned}$$

THEOREM 3.4. Suppose the eigenvalues of P are $\rho_0, \dots, \rho_{n-1}$ with $\rho_0 = 1$, then

$$\lambda_1 \leq \min_{i \neq 0} (1 - R \rho_i).$$

3.2. Definition of Cheeger constants on directed graphs. We have a circulation $F(u, v) = \phi(u) P(u, v)$. Define

$$F(\partial S) = \sum_{u \in S, v \notin S} F(u, v), F(v) = \sum_{u, u \rightarrow v} F(u, v) = \sum_{w, v \rightarrow w} F(v, w), F(S) = \sum_{v \in S} F(v),$$

then $F(\partial S) = F(\partial \bar{S})$.

DEFINITION (Cheeger constant). The Cheeger constant of a graph G is defined as

$$h(G) = \inf_{S \subset V} \frac{F(\partial S)}{\min(F(S), F(\bar{S}))}$$

NOTE. Compare with the undirected graph condition where

$$h_G = \inf_{S \subset V} \frac{|\partial S|}{\min(|S|, |\bar{S}|)}.$$

Similarly, we have

$$\begin{aligned} h_G(\text{undirected}) &= \inf_{S \subset V} \frac{|\partial S|}{\min(|S|, |\bar{S}|)} \\ &= \inf_{S \subset V} \frac{\sum_{u \in S, v \in \bar{S}} w_{uv}}{\min(\sum_{u \in S} d(u), \sum_{u \in \bar{S}} d(u))} \\ h_G(\text{directed}) &= \inf_{S \subset V} \frac{\sum_{u \in S, v \in \bar{S}} \phi(u)P(u, v)}{\min(\sum_{u \in S} \phi(u), \sum_{u \in \bar{S}} \phi(u))} \\ &= \inf_{S \subset V} \frac{F(\partial S)}{\min(F(S), F(\bar{S}))}. \end{aligned}$$

THEOREM 3.5. For every directed graph G ,

$$\frac{h^2(G)}{2} \leq \lambda_1 \leq 2h(G).$$

The proof is similar to the undirected case using Rayleigh quotient and Theorem 3.2.

3.3. An example of the singularity of Cheeger constant on a directed graph. We have already given an example of a directed graph with $n + 1$ vertices and stationary distribution ϕ satisfying $\phi(v_n) = 2^{-n}\phi(v_0)$. Now we make a copy of this graph and denote the new $n + 1$ vertices u_0, \dots, u_n . Joining the two graphs together by two edges $v_n \rightarrow u_n$ and $u_n \rightarrow v_n$, we get a bigger directed graph. Let $S = (v_0, \dots, v_n)$, we have $h(G) \sim 2^{-n}$. In comparison, $h(G) \geq \frac{2}{n(n-1)}$ for undirected graph.

3.4. Estimate the convergence rate of random walks on directed graphs. Define the distance of P after s steps and ϕ as

$$\Delta(s) = \max_{y \in V} \left(\sum_{x \in V} \frac{(P^s(y, x) - \phi(x))^2}{\phi(x)} \right)^{1/2}.$$

Modify the random walk into a lazy random walk $\tilde{P} = \frac{I+P}{2}$, so that it is aperiodic.

THEOREM 3.6.

$$\Delta(t)^2 \leq C(1 - \frac{\lambda_1}{2})^t.$$

3.5. Random Walks on Digraphs, The Generalized Digraph Laplacian, and The Degree of Asymmetry. In this paper the following have been discussed:

- (1) Define an asymmetric Laplacian $\tilde{\mathcal{L}}$ on directed graph;
- (2) Use $\tilde{\mathcal{L}}$ to estimate the hitting time and commute time of the corresponding Markov chain;
- (3) Introduce a metric to measure the asymmetry of $\tilde{\mathcal{L}}$ and use this measure to give a tighter bound on the Markov chain mixing rate and a bound on the Cheeger constant.

Let P be the transition matrix of Markov chain, and $\pi = (\pi_1, \dots, \pi_n)^T$ (column vector) denote its stationary distribution (which is unique if the Markov chain is irreducible, or if the directed graph is strongly connected). Let $\Pi = \text{diag}\{\pi_1, \dots, \pi_n\}$, then we define the normalized Laplacian $\tilde{\mathcal{L}}$ on directed graph:

$$(100) \quad \tilde{\mathcal{L}} = I - \Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}}$$

3.5.1. Hitting time, commute time and fundamental matrix. We establish the relations between $\tilde{\mathcal{L}}$ and the hitting time and commute time of random walk on directed graph through the fundamental matrix $Z = [z_{ij}]$, which is defined as:

$$(101) \quad z_{ij} = \sum_{t=0}^{\infty} (p_{ij}^t - \pi_j), \quad 1 \leq i, j \leq n$$

or alternatively as an infinite sum of matrix series:

$$(102) \quad Z = \sum_{t=0}^{\infty} (P^t - \mathbf{1}\pi^T)$$

With the fundamental matrix, the hitting time and commute time can be expressed as follows:

$$(103) \quad H_{ij} = \frac{z_{jj} - z_{ij}}{\pi_j}$$

$$(104) \quad C_{ij} = H_{ij} + H_{ji} = \frac{z_{jj} - z_{ij}}{\pi_j} + \frac{z_{ii} - z_{ji}}{\pi_i}$$

Using (102), we can write the fundamental matrix Z in a more explicit form. Notice that

$$(105) \quad (P - \mathbf{1}\pi^T)(P - \mathbf{1}\pi^T) = P^2 - \mathbf{1}\pi^T P - P\mathbf{1}\pi^T + \mathbf{1}\pi^T\mathbf{1}\pi^T = P^2 - \mathbf{1}\pi^T$$

We use the fact that $\mathbf{1}$ and π are the right and left eigenvector of the transition matrix P with eigenvalue 1, and that $\pi^T\mathbf{1} = 1$ since π is a distribution. Then

$$(106) \quad Z + \mathbf{1}\pi^T = \sum_{t=0}^{\infty} (P - \mathbf{1}\pi^T)^t = (I - P + \mathbf{1}\pi^T)^{-1}$$

3.5.2. Green's function and Laplacian for directed graph. If we treat the directed graph Laplacian $\tilde{\mathcal{L}}$ as an asymmetric operator on a directed graph G , then we can define the Green's Function $\tilde{\mathcal{G}}$ (without boundary condition) for directed graph. The entries of G satisfy the conditions:

$$(107) \quad (\tilde{\mathcal{G}}\tilde{\mathcal{L}})_{ij} = \delta_{ij} - \sqrt{\pi_i \pi_j}$$

or in the matrix form

$$(108) \quad \tilde{\mathcal{G}}\tilde{\mathcal{L}} = I - \pi^{\frac{1}{2}}\pi^{\frac{1}{2}T}$$

The central theorem in the second paper associate the Green's Function $\tilde{\mathcal{G}}$, the fundamental matrix Z and the normalize directed graph Laplacian $\tilde{\mathcal{L}}$:

THEOREM 3.7. Let $\tilde{\mathcal{Z}} = \Pi^{\frac{1}{2}} Z \Pi^{-\frac{1}{2}}$ and $\tilde{\mathcal{L}}^\dagger$ denote the Moore-Penrose pseudo-inverse $\tilde{\mathcal{L}}$, then

$$(109) \quad \tilde{\mathcal{G}} = \tilde{\mathcal{Z}} = \tilde{\mathcal{L}}^\dagger$$

3.6. measure of asymmetric and its relation to Cheeger constant and mixing rate. To measure the asymmetry in directed graph, we write the $\tilde{\mathcal{L}}$ into the sum of a symmetric part and a skew-symmetric part:

$$(110) \quad \tilde{\mathcal{L}} = \frac{1}{2}(\tilde{\mathcal{L}} + \tilde{\mathcal{L}}^T) + \frac{1}{2}(\tilde{\mathcal{L}} - \tilde{\mathcal{L}}^T)$$

$\frac{1}{2}(\tilde{\mathcal{L}} + \tilde{\mathcal{L}}^T) = \mathcal{L}$ is the symmetrized Laplacian introduced in the first paper. Let $\Delta = \frac{1}{2}(\tilde{\mathcal{L}} - \tilde{\mathcal{L}}^T)$, the Δ captures the difference between $\tilde{\mathcal{L}}$ and its transpose. Let σ_i , λ_i and δ_i ($1 \leq i \leq n$) denotes the i -th singular value of \mathcal{L} , \mathcal{L} , Δ in ascending order ($\sigma_1 = \lambda_1 = \delta_1 = 0$). Then the relation $\tilde{\mathcal{L}} = \mathcal{L} + \Delta$ implies

$$(111) \quad \lambda_i \leq \sigma_i \leq \lambda_i + \delta_n$$

Therefore $\delta_n = \|\Delta\|_2$ is used to measure the degree of asymmetry in the directed graph.

The following two theorems are application of this measure.

THEOREM 3.8. The second singular of $\tilde{\mathcal{L}}$ has bounds :

$$(112) \quad \frac{h(G)^2}{2} \leq \sigma_2 \leq (1 + \frac{\delta_n}{\lambda_2}) \cdot 2h(G)$$

where $h(G)$ is the Cheeger constant of graph G

THEOREM 3.9. For a aperiodic Markov chain P ,

$$(113) \quad \delta_n^2 \leq \max\{\frac{\|\tilde{P}f\|^2}{\|f\|^2} : f \perp \pi^{\frac{1}{2}}\} \leq (1 - \lambda_2)^2 + 2\delta_n\lambda_n + \delta_n^2$$

where $\tilde{P} = \Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}}$

4. Lumpability of Markov Chain

Let P be the transition matrix of a Markov chain on graph $G = (V, E)$ with $V = \{1, 2, \dots, n\}$, i.e. $P_{ij} = \Pr\{x_t = j : x_{t-1} = i\}$. Assume that V admits a partition Ω :

$$\begin{aligned} V &= \bigcup_{i=1}^k \Omega_i, \quad \Omega_i \cap \Omega_j = \emptyset, \quad i \neq j. \\ \Omega &= \{\Omega_s : s = 1, \dots, k\}. \end{aligned}$$

Observe a sequence $\{x_0, x_1, \dots, x_t\}$ sampled from the Markov chain with initial distribution π_0 . Relabel $x_t \mapsto y_t \in \{1, \dots, k\}$ by

$$y_t = \sum_{s=1}^k s \chi_{\Omega_s}(x_t),$$

where χ is the characteristic function. Thus we obtain a sequence (y_t) which is a coarse-grained representation of original sequence.

DEFINITION (Lumpability, Kemeny-Snell 1976). P is lumpable with respect to partition Ω if the sequence $\{y_t\}$ is Markovian. In other words, the transition probabilities do not depend on the choice of initial distribution π_0 and history, i.e.

$$(114) \quad \text{Prob}_{\pi_0}\{x_t \in \Omega_{k_t} : x_{t-1} \in \Omega_{k_{t-1}}, \dots, x_0 \in \Omega_{k_0}\} = \text{Prob}\{x_t \in \Omega_{k_t} : x_{t-1} \in \Omega_{k_{t-1}}\}.$$

The lumpability condition above can be rewritten as

$$(115) \quad \text{Prob}_{\pi_0}\{y_t = k_t : y_{t-1} = k_{t-1}, \dots, y_0 = k_0\} = \text{Prob}\{y_t = k_t : y_{t-1} = k_{t-1}\}.$$

THEOREM 4.1. **I.** (Kemeny-Snell 1976) P is lumpable with respect to partition $\Omega \Leftrightarrow \forall \Omega_s, \Omega_t \in \Omega, \forall i, j \in \Omega_s, \hat{P}_{i\Omega_t} = \hat{P}_{j\Omega_t}$, where $\hat{P}_{i\Omega_t} = \sum_{j \in \Omega_t} P_{ij}$.

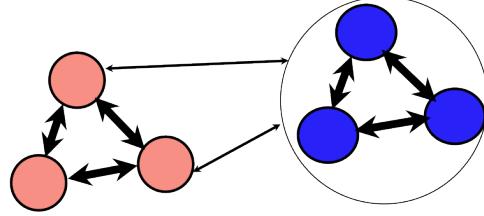


FIGURE 1. Lumpability condition $\hat{P}_{i\Omega_t} = \hat{P}_{j\Omega_t}$

II. (Meila-Shi 2001) P is lumpable with respect to partition Ω and \hat{P} ($\hat{p}_{st} = \sum_{i \in \Omega_s, j \in \Omega_t} p_{ij}$) is nonsingular $\Leftrightarrow P$ has k independent piecewise constant right eigenvectors in $\text{span}\{\chi_{\Omega_s} : s = 1, \dots, k\}$.

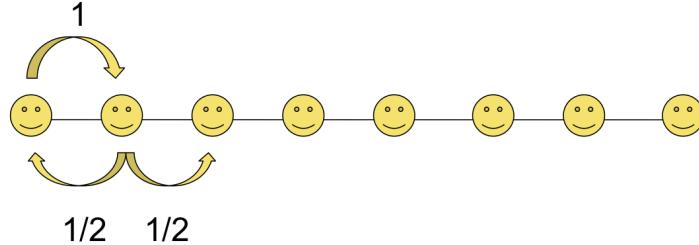


FIGURE 2. A linear chain of $2n$ nodes with a random walk.

EXAMPLE. Consider a linear chain with $2n$ nodes (Figure 2) whose adjacency matrix and degree matrix are given by

$$A = \begin{bmatrix} 0 & 1 & & & \\ 1 & 0 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & 0 & 1 \\ & & & 1 & 0 \end{bmatrix}, \quad D = \text{diag}\{1, 2, \dots, 2, 1\}$$

So the transition matrix is $P = D^{-1}A$ which is illustrated in Figure 2. The spectrum of P includes two eigenvalues of magnitude 1, *i.e.* $\lambda_0 = 1$ and $\lambda_{n-1} = -1$. Although P is not a *primitive* matrix here, it is *lumpable*. Let $\Omega_1 = \{\text{odd nodes}\}$, $\Omega_2 = \{\text{even nodes}\}$. We can check that I and II are satisfied.

To see I, note that for any two even nodes, say $i = 2$ and $j = 4$, $\hat{P}_{i\Omega_2} = \hat{P}_{j\Omega_2} = 1$ as their neighbors are all odd nodes, whence I is satisfied. To see II, note that ϕ_0 (associated with $\lambda_0 = 1$) is a constant vector while ϕ_1 (associated with $\lambda_{n-1} = -1$) is constant on even nodes and odd nodes respectively. Figure 3 shows the lumpable states when $n = 4$ in the left.

Note that lumpable states might not be optimal bi-partitions in $NCUT = Cut(S)/\min(vol(S), vol(\bar{S}))$. In this example, the optimal bi-partition by Ncut is

given by $S = \{1, \dots, n\}$, shown in the right of Figure 3. In fact the second largest eigenvalue $\lambda_1 = 0.9010$ with eigenvector

$$v_1 = [0.4714, 0.4247, 0.2939, 0.1049, -0.1049, -0.2939, -0.4247, -0.4714],$$

give the optimal bi-partition.

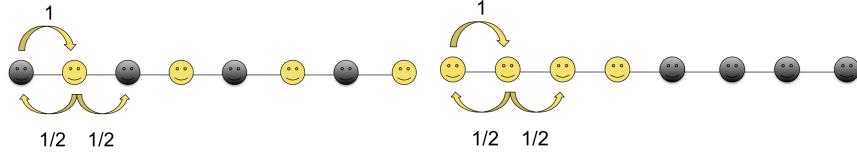


FIGURE 3. Left: two lumpable states; Right: optimal-bipartition of Ncut.

EXAMPLE. Uncoupled Markov chains are lumpable, e.g.

$$P_0 = \begin{bmatrix} \Omega_1 & & \\ & \Omega_2 & \\ & & \Omega_3 \end{bmatrix}, \quad \hat{P}_{it} = \hat{P}_{jt} = 0.$$

A markov chain $\tilde{P} = P_0 + O(\epsilon)$ is called nearly uncoupled Markov chain. Such Markov chains can be approximately represented as uncoupled Markov chains with *metastable states*, $\{\Omega_s\}$, where within metastable state transitions are fast while cross metastable states transitions are slow. Such a separation of scale in dynamics often appears in many phenomena in real lives, such as protein folding, your life transitions *primary schools* \mapsto *middle schools* \mapsto *high schools* \mapsto *college/university* \mapsto *work unit*, etc.

Before the proof of the theorem, we note that condition I is in fact equivalent to

$$(116) \quad VUPV = PV,$$

where U is a k -by- n matrix where each row is a uniform probability that

$$U_{is}^{k \times n} = \frac{1}{|\Omega_s|} \chi_{\Omega_s}(i), \quad i \in V, s \in \Omega,$$

and V is a n -by- k matrix where each column is a characteristic function on Ω_s ,

$$V_{sj}^{n \times k} = \chi_{\Omega_s}(j).$$

With this we have $\hat{P} = UPV$ and $UV = I$. Such a matrix representation will be useful in the derivation of condition II. Now we give the proof of the main theorem.

PROOF. I. “ \Rightarrow ” To see the necessity, P is lumpable w.r.t. partition Ω , then it is necessary that

$$\text{Prob}_{\pi_0}\{x_1 \in \Omega_t : x_0 \in \Omega_s\} = \text{Prob}_{\pi_0}\{y_1 = t : y_0 = s\} = \hat{p}_{st}$$

which does not depend on π_0 . Now assume there are two different initial distribution such that $\pi_0^{(1)}(i) = 1$ and $\pi_0^{(2)}(j) = 1$ for $\forall i, j \in \Omega_s$. Thus

$$\hat{p}_{i\Omega_t} = \text{Prob}_{\pi_0^{(1)}}\{x_1 \in \Omega_t : x_0 \in \Omega_s\} = \hat{p}_{st} = \text{Prob}_{\pi_0^{(2)}}\{x_1 \in \Omega_t : x_0 \in \Omega_s\} = \hat{p}_{j\Omega_t}.$$

“ \Leftarrow ” To show the sufficiency, we are going to show that if the condition is satisfied, then the probability

$$\text{Prob}_{\pi_0}\{y_t = t : y_{t-1} = s, \dots, y_0 = k_0\}$$

depends only on $\Omega_s, \Omega_t \in \Omega$. Probability above can be written as $\text{Prob}_{\pi_{t-1}}(y_t = t)$ where π_{t-1} is a distribution with support only on Ω_s which depends on π_0 and history up to $t-1$. But since $\text{Prob}_i(y_t = t) = \hat{p}_{i\Omega_t} \equiv \hat{p}_{st}$ for all $i \in \Omega_s$, then $\text{Prob}_{\pi_{t-1}}(y_t = t) = \sum_{i \in \Omega_s} \pi_{t-1} \hat{p}_{i\Omega_t} = \hat{p}_{st}$ which only depends on Ω_s and Ω_t .

II.

“ \Rightarrow ”

Since \hat{P} is nonsingular, let $\{\psi_i, i = 1, \dots, k\}$ are independent right eigenvectors of \hat{P} , i.e., $\hat{P}\psi_i = \lambda_i \psi_i$. Define $\phi_i = V\psi_i$, then ϕ_i are independent piecewise constant vectors in $\text{span}\{\chi_{\Omega_i}, i = 1, \dots, k\}$. We have

$$P\phi_i = PV\psi_i = VUPV\psi_i = V\hat{P}\psi_i = \lambda_i V\psi_i = \lambda_i \phi_i,$$

i.e. ϕ_i are right eigenvectors of P .

“ \Leftarrow ”

Let $\{\phi_i, i = 1, \dots, k\}$ be k independent piecewise constant right eigenvectors of P in $\text{span}\{\chi_{\Omega_i}, i = 1, \dots, k\}$. There must be k independent vectors $\psi_i \in \mathbb{R}^k$ that satisfied $\phi_i = V\psi_i$. Then

$$P\phi_i = \lambda_i \phi_i \Rightarrow PV\psi_i = \lambda_i V\psi_i,$$

Multiplying VU to the left on both sides of the equation, we have

$$VUPV\psi_i = \lambda_i VUV\psi_i = \lambda_i V\psi_i = PV\psi_i, \quad (UV = I),$$

which implies

$$(VUPV - PV)\Psi = 0, \quad \Psi = [\psi_1, \dots, \psi_k].$$

Since Ψ is nonsingular due to independence of ψ_i , whence we must have $VUPV = PV$. \square

5. Applications of Lumpability: MNcut and Optimal Reduction of Complex Networks

If the random walk on a graph P has top k nearly piece-wise constant right eigenvectors, then the Markov chain P is approximately lumpable. Some spectral clustering algorithms are proposed in such settings.

5.1. MNcut. Meila-Shi (2001) calls the following algorithm as MNcut, standing for *modified Ncut*. Due to the theory above, perhaps we'd better to call it *multiple spectral clustering*.

- 1) Find top k right eigenvectors $P\Phi_i = \lambda_i \Phi_i, i = 1, \dots, k, \lambda_i = 1 - o(\epsilon)$.
- 2) Embedding $Y^{n \times k} = [\phi_1, \dots, \phi_k] \rightarrow$ diffusion map when $\lambda_i \approx 1$.
- 3) k -means (or other suitable clustering methods) on Y to k -clusters.

5.2. Optimal Reduction and Complex Network.

5.2.1. Random Walk on Graph. Let $G = G(S, E)$ denotes an undirected graph. Here S has the meaning of "states". $|S| = n \gg 1$. Let $A = e(x, y)$ denotes its adjacency matrix, that is,

$$e(x, y) = \begin{cases} 1 & x \sim y \\ 0 & \text{otherwise} \end{cases}$$

Here $x \sim y$ means $(x, y) \in E$. Here, weights on different edges are the same 1. They may be different in some cases.

Now we define a random walk on G . Let

$$p(x, y) = \frac{e(x, y)}{d(x)} \quad \text{where } d(x) = \sum_{y \in S} e(x, y)$$

We can check that $P = p(x, y)$ is a stochastic matrix and (S, P) is a Markov chain. If G is connected, this Markov chain is irreducible and if G is not a tree, the chain is even primitive. We assume G is connected from now on. If it is not, we can focus on each of its connected component. So the Markov chain has unique invariant distribution μ by irreducibility:

$$\mu(x) = \frac{d(x)}{\sum_{z \in S} d(z)} \quad \forall x \in S$$

A Markov chain defined as above is reversible. That is, detailed balance condition is satisfied:

$$\mu(x)p(x, y) = \mu(y)p(y, x) \quad \forall x, y \in S$$

Define an inner product on space \mathcal{L}_μ^2 :

$$\langle f, g \rangle_\mu = \sum_{x \in S} \sum_{y \in S} f(x)g(x)\mu(x) \quad f, g \in \mathcal{L}_\mu^2$$

\mathcal{L}_μ^2 is a Hilbert space with this inner product. If we define an operator T on it:

$$Tf(x) = \sum_{y \in S} p(x, y)f(y) = \mathbb{E}_{[y|x]} f(y)$$

We can check that T is a self adjoint operator on \mathcal{L}_μ^2 :

$$\begin{aligned} \langle Tf(x), g(x) \rangle_\mu &= \sum_{x \in S} Tf(x)g(x)\mu(x) \\ &= \sum_{x \in S} \sum_{y \in S} p(x, y)f(y)g(x)\mu(x) \quad \text{with detailed balance condition} \\ &= \sum_{y \in S} \sum_{x \in S} p(y, x)f(y)g(x)\mu(y) \\ &= \sum_{y \in S} f(y)Tg(y)\mu(y) \\ &= \langle f(x), Tg(x) \rangle_\mu \end{aligned}$$

That means T is self-adjoint. So there is a set of orthonormal basis $\{\phi_j(x)\}_{j=0}^{n-1}$ and a set of eigenvalue $\{\lambda_j\}_{j=0}^{n-1} \subset [-1, 1]$, $1 = \lambda_0 > \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{n-1}$, s.t. $\text{Prob} \phi_j = \lambda_j \phi_j$, $j = 0, 1, \dots, n-1$, and $\langle \phi_i, \phi_j \rangle_\mu = \delta_{ij}$, $\forall i, j = 0, 1, \dots, n-1$. So $\phi_j(x)$ is right

eigenvectors. The corresponding left eigenvectors are denoted by $\{\psi_j(x)\}_{j=0}^{n-1}$. One can obtain that $\psi_j(x) = \phi_j(x)\mu(x)$. In fact, because $T\phi_j = \lambda_j\phi_j$,

$$\begin{aligned}\mu(x) \sum_{y \in S} p(x, y)\phi_j(y) &= \lambda_j\phi_j(x)\mu(x) \quad \text{with detailed balance condition} \\ \sum_{y \in S} p(y, x)\mu(y)\phi_j(y) &= \lambda_j\phi_j(x)\mu(x) \quad \text{that is} \\ \psi_j \text{Prob}(x) &= \sum_{y \in S} p(y, x)\phi(y) = \lambda_j(x)\psi(x)\end{aligned}$$

Generally, T has spectral decomposition

$$p(x, y) = \sum_{i=0}^{n-1} \lambda_i \psi_i(x)\phi(y) = \sum_{i=0}^{n-1} p(x, y)\phi_i(x)\phi_i(y)\mu(x)$$

Since P is a stochastic matrix, we have $\lambda_0 = 1$, the corresponding right eigenvector is $\phi_0(x) \equiv 1$, and left eigenvector is the invariant distribution $\psi_0(x) = \mu(x)$

5.2.2. Optimal Reduction. This section is by [ELVE08]. Suppose the number of states n is very large. The scale of Markov chain is so big that we want a smaller chain to present its behavior. That is, we want to decompose the state space S : Let $S = \bigcup_{i=1}^N S_i$, s.t. $N \ll n$, $S_i \cap S_j = \emptyset, \forall i \neq j$, and define a transition probability \hat{P} on it. We want the Markov chain $(\{S_i\}, \hat{P})$ has similar property as chain (S, P) .

We call $\{S_i\}$ coarse space. The first difficult we're facing is whether $(\{S_i\}, \hat{P})$ really Markovian. We want

$$\Pr(X_{i_{t+1}} \in S_{i_{t+1}} | x_{i_t} \in S_{i_t}, \dots X_0 \in S_{i_0}) = \Pr(X_{i_{t+1}} \in S_{i_{t+1}} | x_{i_t} \in S_{i_t})$$

and this probability is independent of initial distribution. This property is so-called lumpability, which you can refer Lecture 9. Unfortunately, lumpability is a strict constraint that it seldom holds.

So we must modify our strategy of reduction. One choice is to do a optimization with some norm on \mathcal{L}_μ^2 . First, Let us introduce Hilbert-Schmidt norm on \mathcal{L}_μ^2 . Suppose F is an operator on \mathcal{L}_μ^2 , and $Ff(x) = \sum_{y \in S} K(x, y)f(y)\mu(y)$. Here K is called a kernel function. If K is symmetric, F is self adjoint. In fact,

$$\begin{aligned}\langle Ff(x), g(x) \rangle_\mu &= \sum_{x \in S} \sum_{y \in S} K(x, y)f(y)\mu(y)g(x)\mu(x) \\ &= \sum_{y \in S} \sum_{x \in S} K(y, x)f(y)\mu(y)g(x)\mu(x) \\ &= \langle f(x), Fg(x) \rangle_\mu\end{aligned}$$

So F guarantee a spectral decomposition. Let $\{\lambda_j\}_{j=0}^{n-1}$ denote its eigenvalue and $\{\phi_j(x)\}_{j=0}^{n-1}$ denote its eigenvector, then $k(x, y)$ can be represented as $K(x, y) = \sum_{j=0}^{n-1} \lambda_j \phi_j(x)\phi_j(y)$. Hilbert-Schmidt norm of F is defined as follow:

$$\|F\|_{HS}^2 = \text{tr}(F^*F) = \text{tr}(F^2) = \sum_{i=0}^{n-1} \lambda_i^2$$

One can check that $\|F\|_{HS}^2 = \sum_{x,y \in S} K^2(x,y)\mu(x)\mu(y)$. In fact,

$$\begin{aligned} RHS &= \sum_{x,y \in S} \left(\sum_{j=0}^{n-1} \lambda_j \phi_j(x) \phi_j(y) \right)^2 \mu(x)\mu(y) \\ &= \sum_{j=0}^{n-1} \sum_{k=0}^{n-1} \lambda_j \lambda_k \sum_{x,y \in S} \phi_j(x) \phi_k(x) \phi_j(y) \phi_k(y) \mu(x)\mu(y) \\ &= \sum_{j=0}^{n-1} \lambda_j^2 \end{aligned}$$

the last equal sign dues do the orthogonality of eigenvectors. It is clear that if $\mathcal{L}_\mu^2 = \mathcal{L}^2$, Hilbert-Schmidt norm is just Frobenius norm.

Now we can write our T as

$$Tf(x) = \sum_{y \in S} p(x,y)f(y) = \sum_{y \in S} \frac{p(x,y)}{\mu(y)} f(y)\mu(y)$$

and take $K(x,y) = \frac{p(x,y)}{\mu(y)}$. By detailed balance condition, K is symmetric. So

$$\|T\|_{HS}^2 = \sum_{x,y \in S} \frac{p^2(x,y)}{\mu^2(y)} \mu(x)\mu(y) = \sum_{x,y \in S} \frac{\mu(x)}{\mu(y)} p^2(x,y)$$

We'll rename $\|P\|_{HS}$ to $\|P\|_\mu$ in the following paragraphs.

Now go back to our reduction problem. Suppose we have a coarse space $\{S_i\}_{i=1}^N$, and a transition probability $\hat{P}(k,l)$, $k, l = 1, 2, \dots, N$ on it. If we want to compare $(\{S_i\}, \hat{P})$ with (S, P) , we must "lift" the coarse process to fine space. One nature consideration is as follow: if $x \in S_k, y \in S_l$, first, we transit from x to S_l follow the rule $\hat{P}(k,l)$, and in S_l , we transit to y "randomly". To make "randomly" rigorously, one may choose the lifted transition probably as follow:

$$\tilde{P}(x,y) = \sum_{k,l=1}^N 1_{S_k}(x) \hat{P}(k,l) 1_{S_l}(y) \frac{1}{|S_l|}$$

One can check that this \tilde{P} is a stochastic matrix, but it is not reversible. One more convenient choice is transit "randomly" by invariant distribution:

$$\tilde{P}(x,y) = \sum_{k,l=1}^N 1_{S_k}(x) \hat{P}(k,l) 1_{S_l}(y) \frac{\mu(y)}{\hat{\mu}(S_l)}$$

where

$$\hat{\mu}(S_l) = \sum_{z \in S_l} \mu(z)$$

Then you can check this matrix is not only a stochastic matrix, but detailed balance condition also hold provides \hat{P} on $\{S_i\}$ is reversible.

Now let us do some summary. Given a decomposition of state space $S = \bigcup_{i=1}^N S_i$, and a transition probability \hat{P} on coarse space, we may obtain a lifted

transition probability \tilde{P} on fine space. Now we can compare $(\{S_i\}, \hat{P})$ and (S, P) in a clear way: $\|P - \tilde{P}\|_\mu$. So our optimization problem can be defined clearly:

$$E = \min_{S_1 \dots S_N} \min_{\hat{P}} \|P - \hat{P}\|_\mu^2$$

That is, given a partition of S , find the optimal \hat{P} to minimize $\|P - \hat{P}\|_\mu^2$, and find the optimal partition to minimize E .

5.2.3. Community structure of complex network. Given a partition $S = \bigcup_{k=1}^N S_k$, the solution of optimization problem

$$\min_{\hat{p}} \|p - \hat{p}\|_\mu^2$$

is

$$\hat{p}_{kl}^* = \frac{1}{\hat{\mu}(S_k)} \sum_{x \in S_k, y \in S_l} \mu(x)p(x, y)$$

It is easy to show that $\{\hat{p}_{kl}^*\}$ form a transition probability matrix with detailed balance condition:

$$\begin{aligned} \hat{p}_{kl}^* &\geq 0 \\ \sum_l \hat{p}_{kl}^* &= \frac{1}{\hat{\mu}(S_k)} \sum_{x \in S_k} \mu(x) \sum_l \sum_{y \in S_l} p(x, y) \\ &= \frac{1}{\hat{\mu}(S_k)} \sum_{x \in S_k} \mu(x) = 1 \\ \hat{\mu}(S_k) \hat{p}_{kl}^* &= \sum_{x \in S_k, y \in S_l} \mu(x)p(x, y) \\ &= \sum_{x \in S_k, y \in S_l} \mu(y)p(y, x) \\ &= \hat{\mu}(S_l) \hat{p}_{lk}^* \end{aligned}$$

The last equality implies that $\hat{\mu}$ is the invariant distribution of the reduced Markov chain. Thus we find the optimal transition probability in the coarse space. \hat{p}^* has the following property

$$\|p - p^*\|_\mu^2 = \|p\|_\mu^2 - \|\hat{p}^*\|_{\hat{\mu}}^2$$

However, the partition of the original graph is not given in advance, so we need to minimize E^* with respect to all possible partitions. This is a combinatorial optimization problem, which is extremely difficult to find the exact solution. An effective approach to obtain an approximate solution, which inherits ideas of K-means clustering, is proposed as following: First we rewrite E^* as

$$\begin{aligned} E^* &= \sum_{x, y \in S} \frac{\mu(x)}{\mu(y)} |p(x, y) - \sum_{k, l=1}^N 1_{S_k}(x) \frac{\hat{p}_{kl}^*}{\hat{\mu}(S_k)} 1_{S_l}(y) \mu(y)|^2 \\ &= \sum_{k, l=1}^N \sum_{x \in S_k, y \in S_l} \mu(x) \mu(y) \left| \frac{p(x, y)}{\mu(y)} - \frac{\hat{p}_{kl}^*}{\hat{\mu}(S_k)} \right|^2 \\ &\triangleq \sum_{k=1}^N \sum_{x \in S_k} E^*(x, S_k) \end{aligned}$$

where

$$E^*(x, S_k) = \sum_{l=1}^N \sum_{y \in S_l} \mu(x)\mu(y) \left| \frac{p(x, y)}{\mu(y)} - \frac{\hat{p}_{kl}^*}{\hat{\mu}(S_k)} \right|^2$$

Based on above expression, a variation of K-means is designed:

E step: Fix partition $\bigcup_{k=1}^N S_k$, compute \hat{p}^* .

M step: Put x in $S_k^{(n+1)}$ such that

$$E^*(x, S_k) = \min_j E^*(x, S_j)$$

5.2.4. *Extensions: Fuzzy Partition.* This part is in [LLE09, LL11]. It is unnecessary to require that each vertex belong to a definite class. We introduce $\rho_k(x)$ as the probability of a vertex x belonging to class k , and we lift the Markov chain in coarse space to fine space using the following transition probability

$$\tilde{p}(x, y) = \sum_{k,l=1}^N \rho_k(x) \hat{p}_{kl} \rho_l(y) \frac{\mu(y)}{\hat{\mu}_l}$$

Now we solve

$$\min_{\hat{p}} \|p - \tilde{p}\|_\mu^2$$

to obtain a optimal reduction.

5.2.5. *Model selection.* Note the number of partition N should also not be given in advance. But in strategies similar to K-means, the value of minimal E^* is monotone decreasing with N . This means larger N is always preferred.

A possible approach is to introduce another quantity which is monotone increasing with N . We take K-means clustering for example. In K-means clustering, only compactness is reflected. If another quantity indicates separation of centers of each cluster, we can minimize the ratio of compactness and separation to find an optimal N .

6. Mean First Passage Time

Consider a Markov chain P on graph $G = (V, E)$. In this section we study the *mean first passage time* between vertices, which exploits the unnormalized graph Laplacian and will be useful for commute time map against diffusion map.

Definition.

- (1) *First passage time (or hitting time):* $\tau_{ij} := \inf(t \geq 0 | x_t = j, x_0 = i)$;
- (2) *Mean First Passage Time:* $T_{ij} = \mathbb{E}_i \tau_{ij}$;
- (3) $\tau_{ij}^+ := \inf(t > 0 | x_t = j, x_0 = i)$, where τ_{ii}^+ is also called *first return time*;
- (4) $T_{ij}^+ = \mathbb{E}_i \tau_{ij}^+$, where T_{ii}^+ is also called *mean first return time*.

Here \mathbb{E}_i denotes the conditional expectation with fixed initial condition $x_0 = i$.

THEOREM 6.1. Assume that P is irreducible. Let $L = D - W$ be the unnormalized graph Laplacian with Moore-Penrose inverse L^\dagger , where $D = \text{diag}(d_i)$ with $d_i = \sum_{j:j \neq i} W_{ij}$ being the degree of node i . Then

- (1) Mean First Passage Time is given by

$$T_{ii} = 0,$$

$$T_{ij} = \sum_k L_{ik}^\dagger d_k - L_{ij}^\dagger \text{vol}(G) + L_{jj}^\dagger \text{vol}(G) - \sum_k L_{jk}^\dagger d_k, \quad i \neq j.$$

(2) Mean First Return Time is given by

$$T_{ii}^+ = \frac{1}{\pi_i}, \quad T_{ij}^+ = T_{ij}.$$

PROOF. Since P is irreducible, then the stationary distribution is unique, denoted by π . By definition, we have

$$(117) \quad T_{ij}^+ = P_{ij} \cdot 1 + \sum_{k \neq j} P_{ik} (T_{kj}^+ + 1)$$

Let $E = 1 \cdot 1^T$ where $1 \in \mathbb{R}^n$ is a vector with all elements one, $T_d^+ = \text{diag}(T_{ii}^+)$. Then 156 becomes

$$(118) \quad T^+ = E + P(T^+ - T_d^+).$$

For the unique stationary distribution π , $\pi^T P = P$, whence we have

$$\begin{aligned} \pi^T T^+ &= \pi^T 1 \cdot 1^T + \pi^T P(T^+ - T_d^+) \\ \pi^T T^+ &= 1^T + \pi^T T^+ - \pi^T T_d^+ \\ 1 &= T_d^+ \pi \\ T_{ii}^+ &= \frac{1}{\pi_i} \end{aligned}$$

Before proceeding to solve equation (156), we first show its solution is unique.

LEMMA 6.2. P is irreducible $\Rightarrow T^+$ and T are both unique.

PROOF. Assume S is also a solution of equation (157), then

$$\begin{aligned} (I - P)S &= E - P\text{diag}(1/\pi_i) = (I - P)T^+ \\ \Leftrightarrow ((I - P)(T^+ - S)) &= 0. \end{aligned}$$

Therefore for irreducible P , S and T^+ must satisfy

$$\begin{cases} \text{diag}(T^+ - S) = 0 \\ T^+ - S = 1u^T, \quad \forall u \end{cases}$$

which implies $T^+ = S$. T 's uniqueness follows from $T = T^+ - T_d^+$. \square

Now we continue with the proof of the main theorem. Since $T = T^+ - T_d^+$, then (156) becomes

$$\begin{aligned} T &= E + PT - T_d^+ \\ (I - P)T &= E - T_d^+ \\ (I - D^{-1}W)T &= F \\ (D - W)T &= DF \\ LT &= DF \end{aligned}$$

where $F = E - T_d^+$ and $L = D - W$ is the (unnormalized) graph Laplacian. Since L is symmetric and irreducible, we have $L = \sum_{k=1}^n \mu_k \nu_k \nu_k^T$, where $0 = \mu_1 < \mu_2 \leq \dots \leq \mu_n$, $\nu_1 = 1/\|1\|$, $\nu_k^T \nu_l = \delta_{kl}$. Let $L^\dagger = \sum_{k=2}^n \frac{1}{\mu_k} \nu_k \nu_k^T$, L^\dagger is called the *pseudo-inverse* (or *Moore-Penrose inverse*) of L . We can test and verify L^\dagger satisfies the

following four conditions

$$\begin{cases} L^\dagger LL^\dagger = L^\dagger \\ LL^\dagger L = L \\ (LL^\dagger)^T = LL^\dagger \\ (L^\dagger L)^T = L^\dagger L \end{cases}$$

From $LT = D(E - T_d^+)$, multiplying both sides by L^\dagger leads to

$$T = L^\dagger DE - L^\dagger DT_d^+ + 1 \cdot u^T,$$

as $1 \cdot u^T \in \ker(L)$, whence

$$\begin{aligned} T_{ij} &= \sum_{k=1}^n L_{ik}^\dagger d_k - L_{ij}^\dagger d_j \cdot \frac{1}{\pi_j} + u_j \\ u_i &= - \sum_{k=1}^n L_{ik}^\dagger d_k + L_{ii}^\dagger \text{vol}(G), \quad j = i \\ T_{ij} &= \sum_k L_{ik}^\dagger d_k - L_{ij}^\dagger \text{vol}(G) + L_{jj}^\dagger \text{vol}(G) - \sum_k L_{jk}^\dagger d_k \end{aligned}$$

Note that $\text{vol}(G) = \sum_i d_i$ and $\pi_i = d_i/\text{vol}(G)$ for all i . □

As L^\dagger is a positive definite matrix, this leads to the following corollary.

COROLLARY 6.3.

$$(119) \quad T_{ij} + T_{ji} = \text{vol}(G)(L_{ii}^\dagger + L_{jj}^\dagger - 2L_{ij}^\dagger).$$

Therefore the average commute time between i and j leads to an Euclidean distance metric

$$d_c(x_i, x_j) := \sqrt{T_{ij} + T_{ji}}$$

often called *commute time distance*.

7. Transition Path Theory

The transition path theory was originally introduced in the context of continuous-time Markov process on continuous state space [EVE06] and discrete state space [MSVE09], see [EVE10] for a review. Another description of discrete transition path theory for molecular dynamics can be also found in [NSVE⁺09]. The following material is adapted to the setting of discrete time Markov chain with transition probability matrix P [?]. We assume reversibility in the following presentation, which can be extended to non-reversible Markov chains.

Assume that an irreducible Markov Chain on graph $G = (V, E)$ admits the following decomposition $P = D^{-1}W = \begin{pmatrix} P_{ll} & P_{lu} \\ P_{ul} & P_{uu} \end{pmatrix}$. Here $V_l = V_0 \cup V_1$ denotes the labeled vertices with source set V_0 (e.g. reaction state in chemistry) and sink set V_1 (e.g. product state in chemistry), and V_u is the unlabeled vertex set (intermediate states). That is,

- $V_0 = \{i \in V_l : f_i = f(x_i) = 0\}$
- $V_1 = \{i \in V_l : f_i = f(x_i) = 1\}$
- $V = V_0 \cup V_1 \cup V_u$ where $V_l = V_0 \cup V_1$

Given two sets V_0 and V_1 in the state space V , the transition path theory tells how these transitions between the two sets happen (mechanism, rates, etc.). If we view V_0 as a reactant state and V_1 as a product state, then one transition from V_0 to V_1 is a reaction event. The reactive trajectories are those part of the equilibrium trajectory that the system is going from V_0 to V_1 .

Let the hitting time of V_l be

$$\tau_i^k = \inf\{t \geq 0 : x(0) = i, x(t) \in V_k\}, \quad k = 0, 1.$$

The central object in transition path theory is the committor function. Its value at $i \in V_u$ gives the probability that a trajectory starting from i will hit the set V_1 first than V_0 , *i.e.*, the success rate of the transition at i .

PROPOSITION 7.1. For $\forall i \in V_u$, define the *committor function*

$$q_i := \text{Prob}(\tau_i^1 < \tau_i^0) = \text{Prob}(\text{trajectory starting from } x_i \text{ hit } V_1 \text{ before } V_0)$$

which satisfies the following Laplacian equation with Dirichlet boundary conditions

$$(Lq)(i) = [(I - P)q](i) = 0, \quad i \in V_u \\ q_{i \in V_0} = 0, q_{i \in V_1} = 1.$$

The solution is

$$q_u = (D_u - W_{uu})^{-1}W_{ul}q_l.$$

PROOF. By definition,

$$q_i = \text{Prob}(\tau_i^1 < \tau_i^0) = \begin{cases} 1 & x_i \in V_1 \\ 0 & x_i \in V_0 \\ \sum_{j \in V_u} P_{ij}q_j & i \in V_u \end{cases}$$

This is because $\forall i \in V_u$,

$$\begin{aligned} q_i &= \text{Pr}(\tau_{iV_1} < \tau_{iV_0}) \\ &= \sum_j P_{ij}q_j \\ &= \sum_{j \in V_1} P_{ij}q_j + \sum_{j \in V_0} P_{ij}q_j + \sum_{j \in V_u} P_{ij}q_j \\ &= \sum_{j \in V_1} P_{ij} + \sum_{j \in V_u} P_{ij}q_j \end{aligned}$$

$$\therefore q_u = P_{ul}q_l + P_{uu}q_u = D_u^{-1}W_{ul}q_l + D_u^{-1}W_{uu}q_u$$

multiply D_u to both side and reorganize

$$(D_u - W_{uu})q_u = W_{ul}q_l$$

If $D_u - W_{uu}$ is reversible, we get

$$q_u = (D_u - W_{uu})^{-1}W_{ul}q_l.$$

□

The committor function provides natural decomposition of the graph. If $q(x)$ is less than 0.5, i is more likely to reach V_0 first than V_1 ; so that $\{i \mid q(x) < 0.5\}$ gives the set of points that are more attached to set V_0 .

Once the committor function is given, the statistical properties of the reaction trajectories between V_0 and V_1 can be quantified. We state several propositions

characterizing transition mechanism from V_0 to V_1 . The proof of them is an easy adaptation of [EVE06, MSVE09] and will be omitted.

PROPOSITION 7.2 (Probability distribution of reactive trajectories). The probability distribution of reactive trajectories

$$(120) \quad \pi_R(x) = \mathbb{P}(X_n = x, n \in R)$$

is given by

$$(121) \quad \pi_R(x) = \pi(x)q(x)(1 - q(x)).$$

The distribution π_R gives the equilibrium probability that a reactive trajectory visits x . It provides information about the proportion of time the reactive trajectories spend in state x along the way from V_0 to V_1 .

PROPOSITION 7.3 (Reactive current from V_0 to V_1). The reactive current from A to B , defined by

$$(122) \quad J(xy) = \mathbb{P}(X_n = x, X_{n+1} = y, \{n, n+1\} \subset R),$$

is given by

$$(123) \quad J(xy) = \begin{cases} \pi(x)(1 - q(x))P_{xy}q(y), & x \neq y; \\ 0, & \text{otherwise.} \end{cases}$$

The reactive current $J(xy)$ gives the average rate the reactive trajectories jump from state x to y . From the reactive current, we may define the effective reactive current on an edge and transition current through a node which characterizes the importance of an edge and a node in the transition from A to B , respectively.

DEFINITION. The *effective current* of an edge xy is defined as

$$(124) \quad J^+(xy) = \max(J(xy) - J(yx), 0).$$

The *transition current* through a node $x \in V$ is defined as

$$(125) \quad T(x) = \begin{cases} \sum_{y \in V} J^+(xy), & x \in A \\ \sum_{y \in V} J^+(yx), & x \in B \\ \sum_{y \in V} J^+(xy) = \sum_{y \in V} J^+(yx), & x \notin A \cup B \end{cases}$$

In applications one often examines partial transition current through a node connecting two communities $V^- = \{x : q(x) < 0.5\}$ and $V^+ = \{x : q(x) \geq 0.5\}$, e.g. $\sum_{y \in V^+} J^+(xy)$ for $x \in V^-$, which shows relative importance of the node in bridging communities.

The reaction rate ν , defined as the number of transitions from V_0 to V_1 happened in a unit time interval, can be obtained from adding up the probability current flowing out of the reactant state. This is stated by the next proposition.

PROPOSITION 7.4 (Reaction rate). The reaction rate is given by

$$(126) \quad \nu = \sum_{x \in A, y \in V} J(xy) = \sum_{x \in V, y \in B} J(xy).$$

Finally, the committor functions also give information about the time proportion that an equilibrium trajectory comes from A (the trajectory hits A last rather than B).

PROPOSITION 7.5. The proportion of time that the trajectory comes from A (resp. from B) is given by

$$(127) \quad \rho^A = \sum_{x \in V} \pi(x)q(x), \quad \rho^B = \sum_{x \in V} \pi(x)(1 - q(x)).$$

CHAPTER 7

Diffusion Map

Finding meaningful low-dimensional structures hidden in high-dimensional observations is an fundamental task in high-dimensional statistics. The classical techniques for dimensionality reduction, principal component analysis (PCA) and multi-dimensional scaling (MDS), guaranteed to discover the true structure of data lying on or near a linear subspace of the high-dimensional input space. PCA finds a low-dimensional embedding of the data points that best preserves their variance as measured in the high-dimensional input space. Classical MDS finds an embedding that preserves the interpoint distances, equivalent to PCA when those distances are Euclidean [TdL00]. However, these linear techniques cannot adequately handle complex nonlinear data. Recently more emphasis is put on detecting non-linear features in the data. For example, ISOMAP [TdL00] etc. extends MDS by incorporating the geodesic distances imposed by a weighted graph. It defines the geodesic distance to be the sum of edge weights along the shortest path between two nodes. The top n eigenvectors of the geodesic distance matrix are used to represent the coordinates in the new n -dimensional Euclidean space. Nevertheless, as mentioned in [EST09], in practice robust estimation of geodesic distance on a manifold is an awkward problem that require rather restrictive assumptions on the sampling. Moreover, since the MDS step in the ISOMAP algorithm intends to preserve the geodesic distance between points, it provides a correct embedding if submanifold is isometric to a convex open set of the subspace. If the submanifold is not convex, then there exist a pair of points that can not be joined by a straight line contained in the submanifold. Therefore, their geodesic distance can not be equal to the Euclidean distance. Diffusion maps [CLL⁺05] leverages the relationship between heat diffusion and a random walk (Markov Chain); an analogy is drawn between the diffusion operator on a manifold and a Markov transition matrix operating on functions defined on a weighted graph whose nodes were sampled from the manifold. A diffusion map, which maps coordinates between data and diffusion space, aims to re-organize data according to a new metric. In this class, we will discuss this very metric-diffusion distance and its related properties.

1. Diffusion map and Diffusion Distance

Viewing the data points x_1, x_2, \dots, x_n as the nodes of a weighted undirected graph $G = (V, E_W)$ ($W = (W_{ij})$), where the weight W_{ij} is a measure of the similarity between x_i and x_j . There are many ways to define W_{ij} , such as:

- (1) **Heat kernel.** If x_i and x_j are connected, put:

$$(128) \quad W_{ij}^\varepsilon = e^{\frac{-\|x_i - x_j\|^2}{\varepsilon}}$$

with some positive parameter $\varepsilon \in \mathbb{R}_0^+$.

(2) **Cosine Similarity**

$$(129) \quad W_{ij} = \cos(\angle(x_i, x_j)) = \frac{x_i}{\|x_i\|} \cdot \frac{x_j}{\|x_j\|}$$

(3) **Kullback-Leibler divergence.** Assume x_i and x_j are two nonvanishing probability distribution, i.e. $\sum_k x_i^k = 1$ and $x_i^k > 0$. Define *Kullback-Leibler divergence*

$$D^{(KL)}(x_i||x_j) = \sum_k x_i^{(k)} \log \frac{x_i^{(k)}}{x_j^{(k)}}$$

and its symmetrization $\bar{D} = D^{(KL)}(x_i||x_j) + D^{(KL)}(x_j||x_i)$, which measure a kind of ‘distance’ between distributions; *Jensen-Shannon divergence* as the symmetrization of KL-divergence between one distribution and their average,

$$D^{(JS)}(x_i, x_j) = D^{(KL)}(x_i||(x_i + x_j)/2) + D^{(KL)}(x_j||(x_i + x_j)/2)$$

A similarity kernel can be

$$(130) \quad W_{ij} = -D^{(KL)}(x_i||x_j)$$

or

$$(131) \quad W_{ij} = -D^{(JS)}(x_i, x_j)$$

The similarity functions are widely used in various applications. Sometimes the matrix W is positive semi-definite (psd), that for any vector $x \in \mathbb{R}^n$,

$$(132) \quad x^T W x \geq 0.$$

PSD kernels includes heat kernels, cosine similarity kernels, and JS-divergence kernels. But in many other cases (e.g. KL-divergence kernels), similarity kernels are not necessarily PSD. For a PSD kernel, it can be understood as a generalized covariance function; otherwise, diffusions as random walks on similarity graphs will be helpful to disclose their structures.

Define $A := D^{-1}W$, where $D = \text{diag}(\sum_{j=1}^n W_{ij}) \triangleq \text{diag}(d_1, d_2, \dots, d_n)$ for symmetric $W_{ij} = W_{ji} \geq 0$. So

$$(133) \quad \sum_{j=1}^n A_{ij} = 1 \quad \forall i \in \{1, 2, \dots, n\} \quad (A_{ij} \geq 0)$$

whence A is a row Markov matrix of the following discrete time Markov chain $\{X_t\}_{t \in N}$ satisfying

$$(134) \quad P(X_{t+1} = x_j | X_t = x_i) = A_{ij}.$$

1.1. Spectral Properties of A . We may reach a spectral decomposition of A with the aid of the following symmetric matrix S which is similar to A . Let

$$(135) \quad S := D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

which is symmetric and has an eigenvalue decomposition

$$(136) \quad S = V \Lambda V^T, \quad \text{where } VV^T = I_n, \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$$

So

$$A = D^{-1}W = D^{-1}(D^{\frac{1}{2}}SD^{\frac{1}{2}}) = D^{-\frac{1}{2}}SD^{\frac{1}{2}}$$

which is similar to S , whence sharing the same eigenvalues as S . Moreover

$$(137) \quad A = D^{-\frac{1}{2}}V\Lambda V^T D^{\frac{1}{2}} = \Phi\Lambda\Psi^T$$

where $\Phi = D^{-\frac{1}{2}}V$ and $\Psi = D^{\frac{1}{2}}V$ give right and left eigenvectors of A respectively, $A\Phi = \Phi\Lambda$ and $\Psi^T A = \Lambda\Psi^T$, and satisfy $\Psi^T\Phi = I_n$.

The Markov matrix A satisfies the following properties by Perron-Frobenius Theory.

PROPOSITION 1.1. (1) A has eigenvalues $\lambda(A) \subset [-1, 1]$.

(2) A is irreducible, if and only if $\forall(i, j) \exists t$ s.t. $(A^t)_{ij} > 0 \Leftrightarrow \text{Graph } G = (V, E)$ is connected

(3) A is irreducible $\Rightarrow \lambda_{\max} = 1$

(4) A is primitive, if and only if $\exists t > 0$ s.t. $\forall(i, j) (A^t)_{ij} > 0 \Leftrightarrow \text{Graph } G = (V, E)$ is path- t connected, i.e. any pair of nodes are connected by a path of length no more than t

(5) A is irreducible and $\forall i, A_{ii} > 0 \Rightarrow A$ is primitive

(6) A is primitive $\Rightarrow -1 \notin \lambda(A)$

(7) W_{ij} is induced from the heat kernel, or any positive definite function $\Rightarrow \lambda(A) \geq 0$

PROOF. (1) assume λ and v are the eigenvalue and eigenvector of A , so $Av = \lambda v$. Find j_0 s.t. $|v_{j_0}| \geq |v_j|, \forall j \neq j_0$ where v_j is the j -th entry of v . Then:

$$\lambda v_{j_0} = (Av)_{j_0} = \sum_{j=1}^n A_{j_0j}v_j$$

So:

$$|\lambda||v_{j_0}| = |\sum_{j=1}^n A_{j_0j}v_j| \leq \sum_{j=1}^n A_{j_0j}|v_j| \leq |v_{j_0}|.$$

(7) Let $S = D^{-1/2}WD^{-1/2}$. As W is positive semi-definite, so S has eigenvalues $\lambda(S) \geq 0$. Note that $A = D^{-1/2}SD^{1/2}$, i.e. similar to S , whence A shares the same eigenvalues with S . \square

Sort the eigenvalues $1 = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq -1$. Denote $\Phi = [\phi_1, \dots, \phi_n]$ and $\Psi = [\psi_1, \dots, \psi_n]$. So the primary (first) right and left eigenvectors are

$$\phi_1 = \mathbf{1},$$

$$\psi_1 = \pi$$

as the stationary distribution of the Markov chain, respectively.

1.2. Diffusion Map and Distance. Diffusion map of a point x is defined as the weighted Euclidean embedding via right eigenvectors of Markov matrix A . From the interpretation of the matrix A as a Markov transition probability matrix

$$(138) \quad A_{ij} = Pr\{s(t+1) = x_j | s(t) = x_i\}$$

it follows that

$$(139) \quad A_{ij}^t = Pr\{s(t+1) = x_j | s(0) = x_i\}$$

We refer to the i 'th row of the matrix A^t , denoted $A_{i,*}^t$, as the transition probability of a t -step random walk that starts at x_i . We can express A^t using the decomposition of A . Indeed, from

$$(140) \quad A = \Phi \Lambda \Psi^T$$

with $\Psi^T \Phi = I$, we get

$$(141) \quad A^t = \Phi \Lambda^t \Psi^T.$$

Written in a component-wise way, this is equivalent to

$$(142) \quad A_{ij}^t = \sum_{k=1}^n \lambda_k^t \phi_k(i) \psi_k(j).$$

Therefore Φ and Ψ are right and left eigenvectors of A^t , respectively.

Let the diffusion map $\Phi_t : V \mapsto \mathbb{R}^n$ at scale t be

$$(143) \quad \Phi_t(x_i) := \begin{pmatrix} \lambda_1^t \phi_1(i) \\ \lambda_2^t \phi_2(i) \\ \vdots \\ \lambda_n^t \phi_n(i) \end{pmatrix}$$

The mapping of points onto the diffusion map space spanned the right eigenvectors of the row Markov matrix has a well defined probabilistic meaning in terms of the random walks. Lumpable Markov chains with Piece-wise constant right eigenvectors thus help us understand the behavior of diffusion maps and distances in such cases.

The diffusion distance is defined to be the Euclidean distances between embedded points,

$$(144) \quad d_t(x_i, x_j) := \|\Phi_t(x_i) - \Phi_t(x_j)\|_{\mathbb{R}^n} = \left(\sum_{k=1}^n \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2 \right)^{1/2}.$$

The main intuition to define diffusion distance is to describe “perceptual distances” of points in the same and different clusters. For example Figure 1 shows that points within the same cluster have small diffusion distances while in different clusters have large diffusion distances. This is because the metastability phenomenon of random walk on graphs where each cluster represents a metastable state. The main properties of diffusion distances are as follows.

- Diffusion distances reflect average path length connecting points via random walks.
- Small t represents local random walk, where diffusion distances reflect local geometric structure.
- Large t represents global random walk, where diffusion distances reflect large scale cluster or connected components.

1.3. Examples.

Three examples about diffusion map:

EX1: two circles.

Suppose graph $G : (V, E)$. Matrix W satisfies $w_{ij} > 0$, if and only if $(i, j) \in E$. Choose $k(x, y) = I_{\|x-y\| < \delta}$. In this case,

$$A = \begin{pmatrix} A_1 & 0 \\ 0 & A_2 \end{pmatrix},$$

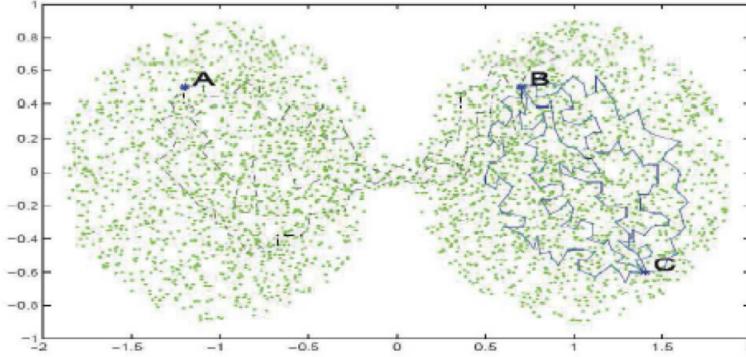


FIGURE 1. Diffusion Distances $d_t(A, B) \gg d_t(B, C)$ while graph shortest path $d_{geod}(A, B) \sim d_{geod}(B, C)$.

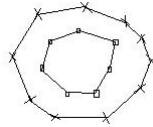


FIGURE 2. Two circles

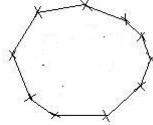


FIGURE 3. EX2 single circle

where A_1 is a $n_1 \times n_1$ matrix, A_2 is a $n_2 \times n_2$ matrix, $n_1 + n_2 = n$.

Notice that the eigenvalue $\lambda_0 = 1$ of A is of multiplicity 2, the two eigenvectors are $\phi_0 = 1_n$ and $\phi_0' = [c_1 1_{n_1}^T, c_2 1_{n_2}^T]^T$ $c_1 \neq c_2$.

$$\text{Diffusion Map : } \begin{cases} \Phi_t^{1D}(x_1), \dots, \Phi_t^{1D}(x_{n_1}) = c_1 \\ \Phi_t^{1D}(x_{n_1+1}), \dots, \Phi_t^{1D}(x_n) = c_2 \end{cases}$$

EX2: ring graph. "single circle"

In this case, W is a circulant matrix

$$W = \begin{pmatrix} 1 & 1 & 0 & 0 & \cdots & 1 \\ 1 & 1 & 1 & 0 & \cdots & 0 \\ 0 & 1 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

The eigenvalue of W is $\lambda_k = \cos \frac{2\pi k}{n}$ $k = 0, 1, \dots, \frac{n}{2}$ and the corresponding eigenvector is $(u_k)_j = e^{i \frac{2\pi}{n} kj}$ $j = 1, \dots, n$. So we can get $\Phi_t^{2D}(x_i) = (\cos \frac{2\pi k j}{n}, \sin \frac{2\pi k j}{n}) c^t$

EX3: order the face. Let



FIGURE 4. Order the face

$$k_\varepsilon(x, y) = \exp\left(-\frac{\|x - y\|^2}{\varepsilon}\right),$$

$W_{ij}^\varepsilon = k_\varepsilon(x_i, x_j)$ and $A_\varepsilon = D^{-1}W^\varepsilon$ where $D = \text{diag}(\sum_j W_{ij}^\varepsilon)$. Define a graph Laplacian (recall that $L = D^{-1}A - I$)

$$L_\varepsilon := \frac{1}{\varepsilon}(A_\varepsilon - I) \xrightarrow{\varepsilon \rightarrow 0} \text{backward Kolmogorov operator}$$

$$L_\varepsilon f = \frac{1}{2} \Delta_M f - \nabla f \cdot \nabla V \Rightarrow L_\varepsilon \phi = \lambda \phi \Rightarrow \begin{cases} \frac{1}{2} \phi''(s) - \phi'(s)V'(s) = \lambda \phi(s) \\ \phi'(0) = \phi'(1) = 0 \end{cases}$$

Where $V(s)$ is the Gibbs free energy and $p(s) = e^{-V(s)}$ is the density of data points along the curve. Δ_M is Laplace-Beltrami Operator. If $p(x) = \text{const}$, we can get

$$(145) \quad V(s) = \text{const} \Rightarrow \phi''(s) = 2\lambda\phi(s) \Rightarrow \phi_k(s) = \cos(k\pi s), 2\lambda_k = -k^2\pi^2$$

On the other hand $p(s) \neq \text{const}$, one can show ¹ that $\phi_1(s)$ is monotonic for arbitrary $p(s)$. As a result, the faces can still be ordered by using $\phi_1(s)$.

1.4. Properties of Diffusion Distance.

LEMMA 1.2. The diffusion distance is equal to a ℓ^2 distance between the probability clouds $A_{i,*}^t$ and $A_{j,*}^t$ with weights $1/d_l$, i.e.,

$$(146) \quad d_t(x_i, x_j) = \|A_{i,*}^t - A_{j,*}^t\|_{\ell^2(\mathbb{R}^n, 1/d)}$$

¹by changing to polar coordinate $p(s)\phi'(s) = r(s)\cos\theta(s)$, $\phi(s) = r(s)\sin\theta(s)$ (the so-called ‘Prufer Transform’) and then try to show that $\phi'(s)$ is never zero on $(0, 1)$.

PROOF.

$$\begin{aligned}
\|A_{i,*}^t - A_{j,*}^t\|_{\ell^2(\mathbb{R}^n, 1/d)}^2 &= \sum_{l=1}^n (A_{il}^t - A_{jl}^t)^2 \frac{1}{d_l} \\
&= \sum_{l=1}^n \left[\sum_{k=1}^n \lambda_k^t \phi_k(i) \psi_k(l) - \lambda_k^t \phi_k(j) \psi_k(l) \right]^2 \frac{1}{d_l} \\
&= \sum_{l=1}^n \sum_{k,k'} \lambda_k^t (\phi_k(i) - \phi_k(j)) \psi_k(l) \lambda_{k'}^t (\phi_{k'}(i) - \phi_{k'}(j)) \psi_{k'}(l) \frac{1}{d_l} \\
&= \sum_{k,k'} \lambda_k^t \lambda_{k'}^t (\phi_k(i) - \phi_k(j)) (\phi_{k'}(i) - \phi_{k'}(j)) \sum_{l=1}^n \frac{\psi_k(l) \psi_{k'}(l)}{d_l} \\
&= \sum_{k,k'} \lambda_k^t \lambda_{k'}^t (\phi_k(i) - \phi_k(j)) (\phi_{k'}(i) - \phi_{k'}(j)) \delta_{kk'} \\
&= \sum_{k=1}^n \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2 \\
&= d_t^2(x_i, x_j)
\end{aligned}$$

□

In practice we usually do not use the mapping Φ_t but rather the truncated diffusion map Φ_t^δ that makes use of fewer than n coordinates. Specifically, Φ_t^δ uses only the eigenvectors for which the eigenvalues satisfy $|\lambda_k|^t > \delta$. When t is enough large, we can use the truncated diffusion distance:

$$(147) \quad d_t^\delta(x_i, x_j) = \|\Phi_t^\delta(x_i) - \Phi_t^\delta(x_j)\| = \left[\sum_{k: |\lambda_k|^t > \delta} \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2 \right]^{\frac{1}{2}}$$

as an approximation of the weighted ℓ^2 distance of the probability clouds. We now derive a simple error bound for this approximation.

LEMMA 1.3 (Truncated Diffusion Distance). The truncated diffusion distance satisfies the following upper and lower bounds.

$$d_t^2(x_i, x_j) - \frac{2\delta^2}{d_{min}} (1 - \delta_{ij}) \leq [d_t^\delta(x_i, x_j)]^2 \leq d_t^2(x_i, x_j),$$

where $d_{min} = \min_{1 \leq i \leq n} d_i$ with $d_i = \sum_j W_{ij}$.

PROOF. Since, $\Phi = D^{-\frac{1}{2}}V$, where V is an orthonormal matrix ($VV^T = V^T V = I$), it follows that

$$(148) \quad \Phi \Phi^T = D^{-\frac{1}{2}} V V^T D^{-\frac{1}{2}} = D^{-1}$$

Therefore,

$$(149) \quad \sum_{k=1}^n \phi_k(i) \phi_k(j) = (\Phi \Phi^T)_{ij} = \frac{\delta_{ij}}{d_i}$$

and

$$(150) \quad \sum_{k=1}^n (\phi_k(i) - \phi_k(j))^2 = \frac{1}{d_i} + \frac{1}{d_j} - \frac{2\delta_{ij}}{d_i}$$

clearly,

$$(151) \quad \sum_{k=1}^n (\phi_k(i) - \phi_k(j))^2 \leq \frac{2}{d_{\min}} (1 - \delta_{ij}), \text{ for all } i, j = 1, 2, \dots, n$$

As a result,

$$\begin{aligned} [d_t^\delta(x_i, x_j)]^2 &= d_t^2(x_i, x_j) - \sum_{k: |\lambda_k|^t < \delta} \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2 \\ &\geq d_t^2(x_i, x_j) - \delta^2 \sum_{k: |\lambda_k|^t < \delta} (\phi_k(i) - \phi_k(j))^2 \\ &\geq d_t^2(x_i, x_j) - \delta^2 \sum_{k=1}^n (\phi_k(i) - \phi_k(j))^2 \\ &\geq d_t^2(x_i, x_j) - \frac{2\delta^2}{d_{\min}} (1 - \delta_{ij}) \end{aligned}$$

on the other hand, it is clear that

$$(152) \quad [d_t^\delta(x_i, x_j)]^2 \leq d_t^2(x_i, x_j)$$

We conclude that

$$(153) \quad d_t^2(x_i, x_j) - \frac{2\delta^2}{d_{\min}} (1 - \delta_{ij}) \leq [d_t^\delta(x_i, x_j)]^2 \leq d_t^2(x_i, x_j)$$

□

Therefore, for small δ the truncated diffusion distance provides a very good approximation to the diffusion distance. Due to the fast decay of the eigenvalues, the number of coordinates used for the truncated diffusion map is usually much smaller than n , especially when t is large.

1.5. Is the diffusion distance really a distance? A distance function $d : X \times X \rightarrow \mathbb{R}$ must satisfy the following properties:

- (1) Symmetry: $d(x, y) = d(y, x)$
- (2) Non-negativity: $d(x, y) \geq 0$
- (3) Identity of indiscernibles: $d(x, y) = 0 \Leftrightarrow x = y$
- (4) Triangle inequality: $d(x, z) + d(z, y) \geq d(x, y)$

Since the diffusion map is an embedding into the Euclidean space \mathbb{R}^n , the diffusion distance inherits all the metric properties of \mathbb{R}^n such as symmetry, non-negativity and the triangle inequality. The only condition that is not immediately implied is $d_t(x, y) = 0 \Leftrightarrow x = y$. Clearly, $x_i = x_j$ implies that $d_t(x_i, x_j) = 0$. But is it true that $d_t(x_i, x_j) = 0$ implies $x_i = x_j$? Suppose $d_t(x_i, x_j) = 0$. Then,

$$(154) \quad 0 = d_t^2(x_i, x_j) = \sum_{k=1}^n \lambda_k^{2t} (\phi_k(i) - \phi_k(j))^2$$

It follows that $\phi_k(i) = \phi_k(j)$ for all k with $\lambda_k \neq 0$. But there is still the possibility that $\phi_k(i) \neq \phi_k(j)$ for k with $\lambda_k = 0$. We claim that this can happen only whenever i and j have the exact same neighbors and proportional weights, that is:

PROPOSITION 1.4. The situation $d_t(x_i, x_j) = 0$ with $x_i \neq x_j$ occurs if and only if node i and j have the exact same neighbors and proportional weights

$$W_{ik} = \alpha W_{jk}, \alpha > 0, \text{ for all } k \in V.$$

PROOF. (Necessity) If $d_t(x_i, x_j) = 0$, then $\sum_{k=1}^n \lambda_k^{2t}(\phi_k(i) - \phi_k(j))^2 = 0$ and $\phi_k(i) = \phi_k(j)$ for k with $\lambda_k \neq 0$. This implies that $d_{t'}(x_i, x_j) = 0$ for all t' , because

$$(155) \quad d_{t'}(x_i, x_j) = \sum_{k=1}^n \lambda_k^{2t'}(\phi_k(i) - \phi_k(j))^2 = 0.$$

In particular, for $t' = 1$, we get $d_1(x_i, x_j) = 0$. But

$$d_1(x_i, x_j) = \|A_{i,*} - A_{j,*}\|_{\ell^2(\mathbb{R}^n, 1/d)},$$

and since $\|\cdot\|_{\ell^2(\mathbb{R}^n, 1/d)}$ is a norm, we must have $A_{i,*} = A_{j,*}$, which implies for each $k \in V$,

$$\frac{W_{ik}}{d_i} = \frac{W_{jk}}{d_j}, \quad \forall k \in V$$

whence $W_{ik} = \alpha W_{jk}$ where $\alpha = d_i/d_j$, as desired.

(Sufficiency) If $A_{i,*} = A_{j,*}$, then $0 = \sum_{k=1}^n (A_{i,k} - A_{j,k})^2/d_k = d_1^2(x_i, x_j) = \sum_{k=1}^n \lambda_k^2(\phi_k(i) - \phi_k(j))^2$ and therefore $\phi_k(i) = \phi_k(j)$ for k with $\lambda_k \neq 0$, from which it follows that $d_t(x_i, x_j) = 0$ for all t . \square

EXAMPLE 11. In a graph with three nodes $V = \{1, 2, 3\}$ and two edges, say $E = \{(1, 2), (2, 3)\}$, the diffusion distance between nodes 1 and 3 is 0. Here the transition matrix is

$$A = \begin{pmatrix} 0 & 1 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 1 & 0 \end{pmatrix}.$$

2. Commute Time Map and Distance

Diffusion distance depends on time scale parameter t which is hard to select in applications. In this section we introduce another closely related distance, namely *commute time distance*, derived from *mean first passage time* between points. For such distances we do not need to choose the time scale t .

Definition.

- (1) *First passage time (or hitting time)*: $\tau_{ij} := \inf(t \geq 0 | x_t = j, x_0 = i)$;
- (2) *Mean First Passage Time*: $T_{ij} = \mathbb{E}_i \tau_{ij}$;
- (3) $\tau_{ij}^+ := \inf(t > 0 | x_t = j, x_0 = i)$, where τ_{ii}^+ is also called *first return time*;
- (4) $T_{ij}^+ = \mathbb{E}_i \tau_{ij}^+$, where T_{ii}^+ is also called *mean first return time*.

Here \mathbb{E}_i denotes the conditional expectation with fixed initial condition $x_0 = i$.

All the below will show that the (average) commute time between x_i and x_j , i.e. $T_{ij} + T_{ji}$, in fact leads to an Euclidean distance metric which can be used for embedding.

THEOREM 2.1. $d_c(x_i, x_j) := \sqrt{T_{ij} + T_{ji}}$ is an Euclidean distance metric, called *commute time distance*.

PROOF. For simplicity, we will assume that P is irreducible such that the stationary distribution is unique. We will give a constructive proof that $T_{ij} + T_{ji}$ is a squared distance of some Euclidean coordinates for x_i and x_j .

By definition, we have

$$(156) \quad T_{ij}^+ = P_{ij} \cdot 1 + \sum_{k \neq j} P_{ik} (T_{kj}^+ + 1)$$

Let $E = 1 \cdot 1^T$ where $1 \in \mathbb{R}^n$ is a vector with all elements one, $T_d^+ = \text{diag}(T_{ii}^+)$. Then 156 becomes

$$(157) \quad T^+ = E + P(T^+ - T_d^+).$$

For the unique stationary distribution π , $\pi^T P = P$, whence we have

$$\begin{aligned} \pi^T T^+ &= \pi^T 1 \cdot 1^T + \pi^T P(T^+ - T_d^+) \\ \pi^T T^+ &= 1^T + \pi^T T^+ - \pi^T T_d^+ \\ 1 &= T_d^+ \pi \\ T_{ii}^+ &= \frac{1}{\pi_i} \end{aligned}$$

Before proceeding to solve equation (156), we first show its solution is unique.

LEMMA 2.2. P is irreducible $\Rightarrow T^+$ and T are both unique.

PROOF. Assume S is also a solution of equation (157), then

$$\begin{aligned} (I - P)S &= E - P\text{diag}(1/\pi_i) = (I - P)T^+ \\ \Leftrightarrow ((I - P)(T^+ - S)) &= 0. \end{aligned}$$

Therefore for irreducible P , S and T^+ must satisfy

$$\begin{cases} \text{diag}(T^+ - S) = 0 \\ T^+ - S = 1u^T, \quad \forall u \end{cases}$$

which implies $T^+ = S$. T 's uniqueness follows from $T = T^+ - T_d^+$. \square

Now we continue with the proof of the main theorem. Since $T = T^+ - T_d^+$, then (156) becomes

$$\begin{aligned} T &= E + PT - T_d^+ \\ (I - P)T &= E - T_d^+ \\ (I - D^{-1}W)T &= F \\ (D - W)T &= DF \\ LT &= DF \end{aligned}$$

where $F = E - T_d^+$ and $L = D - W$ is the (unnormalized) graph Laplacian. Since L is symmetric and irreducible, we have $L = \sum_{k=1}^n \mu_k \nu_k \nu_k^T$, where $0 = \mu_1 < \mu_2 \leq \dots \leq \mu_n$, $\nu_1 = 1/\|1\|$, $\nu_k^T \nu_l = \delta_{kl}$. Let $L^+ = \sum_{k=2}^n \frac{1}{\mu_k} \nu_k \nu_k^T$, L^+ is called the *pseudo-inverse* (or *Moore-Penrose inverse*) of L . We can test and verify L^+ satisfies the following four conditions

$$\begin{cases} L^+ LL^+ = L^+ \\ LL^+ L = L \\ (LL^+)^T = LL^+ \\ (L^+ L)^T = L^+ L \end{cases}$$

From $LT = D(E - T_d^+)$, multiplying both sides by L^+ leads to

$$T = L^+ DE - L^+ DT_d^+ + 1 \cdot u^T,$$

TABLE 1. Comparisons between diffusion map and commute time map. Here $x \sim y$ means that x and y are in the same cluster and $x \not\sim y$ for different clusters.

Diffusion Map	Commute Time Map
P 's right eigenvectors scale parameters: t and ε $\exists t$, s.t. $x \sim y$, $d_t(x, y) \rightarrow 0$ and $x \not\sim y$, $d_t(x, y) \rightarrow \infty$	L^+ 's eigenvectors scale: ε $x \sim y$, $d_c(x, y)$ small and $x \not\sim y$, $d_c(x, y)$ large?

as $1 \cdot u^T \in \ker(L)$, whence

$$\begin{aligned} T_{ij} &= \sum_{k=1}^n L_{ik}^+ d_k - L_{ij}^+ d_j \cdot \frac{1}{\pi_j} + u_j \\ u_i &= - \sum_{k=1}^n L_{ik}^+ d_k + L_{ii}^+ \text{vol}(G), \quad j = i \\ T_{ij} &= \sum_k L_{ik}^+ d_k - L_{ij}^+ \text{vol}(G) + L_{jj}^+ \text{vol}(G) - \sum_k L_{jk}^+ d_k \end{aligned}$$

Note that $\text{vol}(G) = \sum_i d_i$ and $\pi_i = d_i / \text{vol}(G)$ for all i .

Then

$$(158) \quad T_{ij} + T_{ji} = \text{vol}(G)(L_{ii}^+ + L_{jj}^+ - 2L_{ij}^+).$$

To see it is a squared Euclidean distance, we need the following lemma.

LEMMA 2.3. If K is a symmetric and positive semidefinite matrix, then

$$K(x, x) + K(y, y) - 2K(x, y) = d^2(\Phi(x), \Phi(y)) = \langle \Phi(x), \Phi(x) \rangle + \langle \Phi(y), \Phi(y) \rangle - 2\langle \Phi(x), \Phi(y) \rangle$$

where $\Phi = (\phi_i : i = 1, \dots, n)$ are orthonormal eigenvectors with eigenvalues $\mu_i \geq 0$, such that $K(x, y) = \sum_i \mu_i \phi_i(x) \phi_i(y)$.

Clearly L^+ is a positive semidefinite matrix and we define the *commute time map* by its eigenvectors,

$$\Psi(x_i) = \left(\frac{1}{\sqrt{\mu_2}} \nu_2(i), \dots, \frac{1}{\sqrt{\mu_n}} \nu_n(i) \right)^T \in \mathbb{R}^{n-1}.$$

then $L_{ii}^+ + L_{jj}^+ - 2L_{ij}^+ = \|\Psi(x_i) - \Psi(x_j)\|_{l^2}^2$, and we call $d_r(x_i, x_j) = \sqrt{L_{ii}^+ + L_{jj}^+ - 2L_{ij}^+}$ the *resistance distance*.

So we have $d_c(x_i, x_j) = \sqrt{T_{ij} + T_{ji}} = \sqrt{\text{vol}(G)} d_r(x_i, x_j)$. \square

2.1. Comparisons between diffusion map and commute time map.

However, recently Radl, von Luxburg, and Hein give a *negative* answer for the last desired property of $d_c(x, y)$ in geometric random graphs. Their result is as follows. Let $\mathcal{X} \subseteq \mathbb{R}^p$ be a compact set and let $k : \mathcal{X} \times \mathcal{X} \rightarrow (0, +\infty)$ be a symmetric and continuous function. Suppose that $(x_i)_{i \in \mathbb{N}}$ is a sequence of data points drawn i.i.d. from \mathcal{X} according to a density function $p > 0$ on \mathcal{X} . Define $W_{ij} = k(x_i, x_j)$, $P = D^{-1}W$, and $L = D - W$. Then Radl et al. shows

$$\lim_{n \rightarrow \infty} n d_r(x_i, x_j) = \frac{1}{d(x_i)} + \frac{1}{d(x_j)}$$

where $d(x) = \int_X k(x, y) dp(y)$ is a smoothed density at x , $d_r(x_i, x_j) = \frac{d_c(x_i, x_j)}{\sqrt{\text{vol}(G)}}$ is the resistance distance. This result shows that in this setting commute time distance has no information about cluster information about point cloud data, instead it simply reflects density information around the two points.

3. Diffusion Map: Convergence Theory

Diffusion distance depends on both the geometry and density of the dataset. The key concepts in the analysis of these methods, that incorporates the density and geometry of a dataset. This section we will prove the convergence of diffusion map with heat kernels to its geometric limit, the eigenfunctions of Laplacian-Beltrami operators.

This is left by previous lecture. W is positive definite if using Gaussian Kernel.

One can check that, when

$$Q(x) = \int_{\mathbb{R}} e^{-ix\xi} d\mu(\xi),$$

for some positive finite Borel measure $d\mu$ on \mathbb{R} , then the (symmetric/Hermitian) integral kernel

$$k(x, y) = Q(x - y)$$

is positive definite, that is, for any function $\phi(x)$ on \mathbb{R} ,

$$\int \int \bar{\phi}(x)\phi(y)k(x, y) \geq 0.$$

Proof omitted. The reverse is also true, which is Bochner theorem. High dimensional case is similar.

Take 1-dimensional as an example. Since the Gaussian distribution $e^{-\xi^2/2} d\xi$ is a positive finite Borel measure, and the Fourier transform of Gaussian kernel is itself, we know that $k(x, y) = e^{-|x-y|^2/2}$ is a positive definite integral kernel. The matrix W as an discretized version of $k(x, y)$ keeps the positive-definiteness (make this rigorous? Hint: take $\phi(x)$ as a linear combination of n delta functions).

3.1. Main Result. In this lecture, we will study the bias and variance decomposition for sample graph Laplacians and their asymptotic convergence to Laplacian-Beltrami operators on manifolds.

Let \mathcal{M} be a smooth manifold without boundary in \mathbb{R}^p (e.g. a d -dimensional sphere). Randomly draw a set of n data points, $x_1, \dots, x_n \in M \subset \mathbb{R}^p$, according to distribution $p(x)$ in an independent and identically distributed (i.i.d.) way. We can extract an $n \times n$ weight matrix W_{ij} as follows:

$$W_{ij} = k(x_i, x_j)$$

where $k(x, y)$ is a symmetric $k(x, y) = k(y, x)$ and positivity-preserving kernel $k(x, y) \geq 0$. As an example, it can be the *heat kernel* (or Gaussian kernel),

$$k_\epsilon(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\epsilon}\right),$$

where $\|\cdot\|^2$ is the Euclidean distance in space R^p and ϵ is the bandwidth of the kernel. W_{ij} stands for similarity function between x_i and x_j . A diagonal matrix D is defined with diagonal elements are the row sums of W :

$$D_{ii} = \sum_{j=1}^n W_{ij}.$$

Let's consider a family of re-weighted similarity matrix, with superscript (α) ,

$$W^{(\alpha)} = D^{-\alpha} W D^{-\alpha}$$

and

$$D_{ii}^{(\alpha)} = \sum_{j=1}^n W_{ij}^{(\alpha)}.$$

Denote $A^{(\alpha)} = (D^{(\alpha)})^{-1} W$, and we can verify that $\sum_{j=1}^n A_{ij}^{(\alpha)} = 1$, i.e.a row Markov matrix. Now define $L^{(\alpha)} = A^{(\alpha)} - I = (D^{(\alpha)})^{-1} W^{(\alpha)} - I$; and

$$L_{\epsilon,\alpha} = \frac{1}{\epsilon}(A_\epsilon^{(\alpha)} - I)$$

when $k_\epsilon(x, y)$ is used in constructing W . In general, $L^{(\alpha)}$ and $L_{\epsilon,\alpha}$ are both called *graph Laplacians*. In particular $L^{(0)}$ is the unnormalized graph Laplacian in literature.

The target is to show that graph Laplacian $L_{\epsilon,\alpha}$ converges to continuous differential operators acting on smooth functions on \mathcal{M} the manifold. The convergence can be roughly understood as: we say a sequence of n -by- n matrix $L^{(n)}$ as $n \rightarrow \infty$ converges to a limiting operator \mathcal{L} , if for \mathcal{L} 's eigenfunction $f(x)$ (a smooth function on \mathcal{M}) with eigenvalue λ , that is

$$\mathcal{L}f = \lambda f,$$

the length- n vector $f^{(n)} = (f(x_i)), (i = 1, \dots, n)$ is approximately an eigenvector of $L^{(n)}$ with eigenvalue λ , that is

$$L^{(n)} f^{(n)} = \lambda f^{(n)} + o(1),$$

where $o(1)$ goes to zero as $n \rightarrow \infty$.

Specifically, (the convergence is in the sense of multiplying a positive constant)

- (I) $L_{\epsilon,0} = \frac{1}{\epsilon}(A_\epsilon - I) \rightarrow \frac{1}{2}(\Delta_{\mathcal{M}} + 2\frac{\nabla p}{p} \cdot \nabla)$ as $\epsilon \rightarrow 0$ and $n \rightarrow \infty$. $\Delta_{\mathcal{M}}$ is the Laplace-Beltrami operator of manifold M . At a point on M which is d -dimensional, in local (orthogonal) geodesic coordinate s_1, \dots, s_d , the Laplace-Beltrami operator has the same form as the laplace in calculus

$$\Delta_{\mathcal{M}} f = \sum_{i=1}^d \frac{\partial^2}{\partial s_i^2} f;$$

∇ denotes the gradient of a function on M , and \cdot denotes the inner product on tangent spaces of \mathcal{M} . Note that $p = e^{-V}$, so $\frac{\nabla p}{p} = -\nabla V$.

(Ignore this part if you don't know stochastic process) Suppose we have the following diffusion process

$$dX_t = -\nabla V(X_t)dt + \sigma dW_t^{(M)},$$

where $W_t^{(M)}$ is the Brownian motion on M , and σ is the volatility, say a positive constant, then the backward Kolmogorov operator/Fokker-Plank

operator/infinitesimal generator of the process is

$$\frac{\sigma^2}{2} \Delta_{\mathcal{M}} - \nabla V \cdot \nabla,$$

so we say in (I) the limiting operator is the Fokker-Plank operator. Notice that in Lafon '06 paper they differ the case of $\alpha = 0$ and $\alpha = 1/2$, and argue that only in the later case the limiting operator is the Fokker-Plank. However the difference between $\alpha = 0$ and $\alpha = 1/2$ is a $1/2$ factor in front of $-\nabla V$, and that can be unified by changing the volatility σ to another number. (Actually, according to Thm 2. on Page 15 of Lafon'06, one can check that $\sigma^2 = \frac{1}{1-\alpha}$.) So here we say for $\alpha = 0$ the limiting operator is also Fokker-Plank. (not talked in class, open to discussion...)

- (II) $L_{\epsilon,1} = \frac{1}{\epsilon}(A_{\epsilon}^{(1)} - I) \rightarrow \frac{1}{2}\Delta_{\mathcal{M}}$ as $\epsilon \rightarrow 0$ and $n \rightarrow \infty$. Notice that this case is of important application value: whatever the density $p(x)$ is, the Laplacian-Beltrami operator of \mathcal{M} is approximated, so the geometry of the manifold can be understood.

A special case is that samples x_i are uniformly distributed on \mathcal{M} , whence $\nabla p = 0$. Then (I) and (II) are the same up to multiplying a positive constant, due to that D 's diagonal entries are almost the same number and the re-weight does not do anything.

Convergence results like these can be found in Coifman and Lafon [CL06], *Diffusion maps, Applied and Computational Harmonic Analysis*.

We also refer [Sin06] *From graph to manifold Laplacian: The convergence rate, Applied and Computational Harmonic Analysis* for a complete analysis of the variance error, while the analysis of bias is very brief in this paper.

3.2. Proof. For a smooth function $f(x)$ on \mathcal{M} , let $f = (f_i) \in \mathbb{R}^n$ as a vector defined by $f_i = f(x_i)$. At a given fixed point x_i , we have the formula:

$$\begin{aligned} (Lf)^i &= \frac{1}{\epsilon} \left(\frac{\sum_{j=1}^n W_{ij} f_j}{\sum_{j=1}^n W_{ij}} - f_i \right) = \frac{1}{\epsilon} \left(\frac{\frac{1}{n} \sum_{j=1}^n W_{ij} f_j}{\frac{1}{n} \sum_{j=1}^n W_{ij}} - f_i \right) \\ &= \frac{1}{\epsilon} \left(\frac{\frac{1}{n} \sum_{j \neq i} k_{\epsilon}(x_i, x_j) \cdot f(x_j)}{\frac{1}{n} \sum_{j \neq i} k_{\epsilon}(x_i, x_j)} - f(x_i) + f(x_i)O(\frac{1}{n\epsilon^{\frac{d}{2}}}) \right) \end{aligned}$$

where in the last step the diagonal terms $j = i$ are excluded from the sums resulting in an $O(n^{-1}\epsilon^{-\frac{d}{2}})$ error. Later we will see that compared to the variance error, this term is negligible.

We rewrite the Laplacian above as

$$(159) \quad (Lf)^i = \frac{1}{\epsilon} \left(\frac{F(x_i)}{G(x_i)} - f(x_i) + f(x_i)O(\frac{1}{n\epsilon^{\frac{d}{2}}}) \right)$$

where

$$F(x_i) = \frac{1}{n} \sum_{j \neq i} k_{\epsilon}(x_i, x_j) f(x_j), \quad G(x_i) = \frac{1}{n} \sum_{j \neq i} k_{\epsilon}(x_i, x_j).$$

depends only on the other $n - 1$ data points than x_i . In what follows we treat x_i as a fixed chosen point and write as x .

Bias-Variance Decomposition. The points $x_j, j \neq i$ are independent identically distributed (i.i.d), therefore every term in the summation of $F(x)$ ($G(x)$) are i.i.d., and by the Law of Large Numbers (LLN) one should expect $F(x) \approx \mathbb{E}_{x_1}[k(x, x_1)f(x_1)] = \int_{\mathcal{M}} k(x, y)f(y)p(y)dy$ (and $G(x) \approx \mathbb{E}k(x, x_1) = \int_{\mathcal{M}} k(x, y)p(y)dy$). Recall that given a random variable x , and a sample estimator $\hat{\theta}$ (e.g. sample mean), the bias-variance decomposition is given by

$$\mathbb{E}\|x - \hat{\theta}\|^2 = \mathbb{E}\|x - \mathbb{E}x\|^2 + \mathbb{E}\|\mathbb{E}x - \hat{\theta}\|^2.$$

If we use the same strategy here (though not exactly the same, since $\mathbb{E}\left[\frac{F}{G}\right] \neq \frac{\mathbb{E}[F]}{\mathbb{E}[G]}$!), we can decompose Eqn. (159) as

$$(Lf)^i = \frac{1}{\epsilon} \left(\frac{\mathbb{E}[F]}{\mathbb{E}[G]} - f(x_i) + f(x_i)O\left(\frac{1}{n\epsilon^{\frac{d}{2}}}\right) \right) + \frac{1}{\epsilon} \left(\frac{F(x_i)}{G(x_i)} - \frac{\mathbb{E}[F]}{\mathbb{E}[G]} \right)$$

$$= \text{bias} + \text{variance}.$$

In the below we shall show that for case (I) the estimates are

$$(160) \quad \text{bias} = \frac{1}{\epsilon} \left(\frac{\mathbb{E}[F]}{\mathbb{E}[G]} - f(x) + f(x_i)O\left(\frac{1}{n\epsilon^{\frac{d}{2}}}\right) \right) = \frac{m_2}{2} (\Delta_{\mathcal{M}}f + 2\nabla f \cdot \frac{\nabla p}{p}) + O(\epsilon) + O\left(n^{-1}\epsilon^{-\frac{d}{2}}\right).$$

$$(161) \quad \text{variance} = \frac{1}{\epsilon} \left(\frac{F(x_i)}{G(x_i)} - \frac{\mathbb{E}[F]}{\mathbb{E}[G]} \right) = O(n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}-1}),$$

whence

$$\text{bias} + \text{variance} = O(\epsilon, n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}-1}) = C_1\epsilon + C_2n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}-1}.$$

As the bias is a monotone increasing function of ϵ while the variance is decreasing w.r.t. ϵ , the optimal choice of ϵ is to balance the two terms by taking derivative of the right hand side equal to zero (or equivalently setting $\epsilon \sim n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}-1}$) whose solution gives the optimal rates

$$\epsilon^* \sim n^{-1/(2+d/2)}.$$

[CL06] gives the bias and [HAvL05] contains the variance parts, which are further improved by [Sin06] in both bias and variance.

3.3. The Bias Term.

Now focus on $\mathbb{E}[F]$

$$\mathbb{E}[F] = \mathbb{E} \left[\frac{1}{n} \sum_{j \neq i} k_{\epsilon}(x_i, x_j) f(x_j) \right] = \frac{n-1}{n} \int_{\mathcal{M}} k_{\epsilon}(x, y) f(y) p(y) dy$$

$\frac{n-1}{n}$ is close to 1 and is treated as 1.

- (1) the case of one-dimensional and flat (which means the manifold \mathcal{M} is just a real line, i.e. $\mathcal{M} = \mathbb{R}$)

Let $\tilde{f}(y) = f(y)p(y)$, and $k_{\epsilon}(x, y) = \frac{1}{\sqrt{\epsilon}} e^{-\frac{(x-y)^2}{2\epsilon}}$, by change of variable

$$y = x + \sqrt{\epsilon}z,$$

we have

$$\square = \int_{\mathbb{R}} \tilde{f}(x + \sqrt{\epsilon}z) e^{-\frac{z^2}{2}} dz = m_0 \tilde{f}(x) + \frac{1}{2} m_2 f''(x)\epsilon + O(\epsilon^2)$$

where $m_0 = \int_{\mathbb{R}} e^{-\frac{z^2}{2}} dz$, and $m_2 = \int_{\mathbb{R}} z^2 e^{-\frac{z^2}{2}} dz$.

(2) 1 Dimensional & Not flat:

Divide the integral into 2 parts:

$$\int_m k_\epsilon(x, y) \tilde{f}(y) p(y) dy = \int_{||x-y||>c\sqrt{\epsilon}} \cdot + \int_{||x-y||<c\sqrt{\epsilon}} \cdot$$

First part = \circ

$$|\circ| \leq \|\tilde{f}\|_\infty \frac{1}{\epsilon^{\frac{a}{2}}} e^{-\frac{\epsilon^2}{2\epsilon}},$$

due to $||x - y||^2 > c\sqrt{\epsilon}$

$$c \sim \ln(\frac{1}{\epsilon}).$$

so this item is tiny and can be ignored.

Locally, that is $u \sim \sqrt{\epsilon}$, we have the curve in a plane and has the following parametrized equation

$$(x(u), y(u)) = (u, au^2 + qu^3 + \dots),$$

then the chord length

$$\frac{1}{\epsilon} ||x - y||^2 = \frac{1}{\epsilon} [u^2 + (au^2 + qu^3 + \dots)^2] = \frac{1}{\epsilon} [u^2 + a^2 u^4 + q_5(u) + \dots],$$

where we mark $a^2 u^4 + 2a qu^5 + \dots = q_5(u)$. Next, change variable $\frac{u}{\sqrt{\epsilon}} = z$, then with $h(\xi) = e^{-\frac{\xi}{2}}$

$$h(\frac{||x - y||}{\epsilon})^2 = h(z^2) + h'(z^2)(\epsilon^2 az^4 + \epsilon^{\frac{3}{2}} q_5 + O(\epsilon^2)),$$

also

$$\tilde{f}(s) = \tilde{f}(x) + \frac{d\tilde{f}}{ds}(x)s + \frac{1}{2} \frac{d^2\tilde{f}}{ds^2}(x)s^2 + \dots$$

and

$$s = \int_0^u \sqrt{1 + (2au + 3quu^2 + \dots)^2} du + \dots$$

and

$$\frac{ds}{du} = 1 + 2a^2 u^2 + q_2(u) + O(\epsilon^2), \quad s = u + \frac{2}{3} a^2 u^3 + O(\epsilon^2).$$

Now come back to the integral

$$\begin{aligned} & \int_{||x-y||<c\sqrt{\epsilon}} \frac{1}{\sqrt{\epsilon}} h(\frac{x-y}{\epsilon}) \tilde{f}(s) ds \\ & \approx \int_{-\infty}^{+\infty} [h(z^2) + h'(z^2)(\epsilon^2 az^4 + \epsilon^{\frac{3}{2}} q_5)] \cdot [\tilde{f}(x) + \frac{d\tilde{f}}{ds}(x)(\sqrt{\epsilon}z + \frac{2}{3} a^2 z^2 \epsilon^{\frac{3}{2}}) \\ & \quad + \frac{1}{2} \frac{d^2\tilde{f}}{ds^2}(x)\epsilon z^2] \cdot [1 + 2a^2 + \epsilon^3 y_3(z)] dz \\ & = m_0 \tilde{f}(x) + \epsilon \frac{m_2}{2} (\frac{d^2\tilde{f}}{ds^2}(x) + a^2 \tilde{f}(x)) + O(\epsilon^2), \end{aligned}$$

where the $O(\epsilon^2)$ tails are omitted in middle steps, and $m_0 = \int h(z^2) dz, m_2 = \int z^2 h(z^2) dz$, are positive constants. In what follows we normalize both of

them by m_0 , so only m_2 appears as coefficient in the $O(\epsilon)$ term. Also the fact that $h(\xi) = e^{-\frac{\xi}{2}}$, and so $h'(\xi) = -\frac{1}{2}h(\xi)$, is used.

- (3) For high dimension, \mathcal{M} is of dimension d ,

$$k_\epsilon(x, y) = \frac{1}{\epsilon^{\frac{d}{2}}} e^{-\frac{|x-y|^2}{2\epsilon}},$$

the corresponding result is (Lemma 8 in Appendix B of Lafon '06 paper)

$$(162) \quad \int_{\mathcal{M}} k_\epsilon(x, y) \tilde{f}(y) dy = \tilde{f}(x) + \epsilon \frac{m_2}{2} (\Delta_{\mathcal{M}} \tilde{f} + E(x) \tilde{f}(x)) + O(\epsilon^2),$$

where

$$E(x) = \sum_{i=1}^d a_i(x)^2 - \sum_{i_1 \neq i_2} a_{i_1}(x) a_{i_2}(x),$$

and $a_i(x)$ are the curvatures along coordinates s_i ($i = 1, \dots, d$) at point x .

Now we study the limiting operator and the bias error:

$$\begin{aligned} \frac{\mathbb{E}F}{\mathbb{E}G} &= \frac{\int k_\epsilon(x, y) f(y) p(y) dy}{\int k_\epsilon(x, y) p(y) dy} \approx \frac{f + \epsilon \frac{m_2}{2} (f'' + 2f' \frac{p'}{p} + f \frac{p^2}{p} + Ef) + O(\epsilon^2)}{1 + \epsilon \frac{m_2}{2} (\frac{p''}{p} + E) + O(\epsilon^2)} \\ (163) \quad &= f(x) + \epsilon \frac{m_2}{2} (f'' + 2f' \frac{p'}{p}) + o(\epsilon^2), \end{aligned}$$

and as a result, for generally d -dim case,

$$\frac{1}{\epsilon} \left(\frac{\mathbb{E}F}{\mathbb{E}G} - f(x) \right) = \frac{m_2}{2} (\Delta_{\mathcal{M}} f + 2\nabla f \cdot \frac{\nabla p}{p}) + O(\epsilon).$$

Using the same method and use Eqn. (162), one can show that for case (II) where $\alpha = 1$, the limiting operator is exactly the Laplace-Beltrami operator and the bias error is again $O(\epsilon)$ (homework).

About \mathcal{M} with boundary: firstly the limiting differential operator bears Newmann/no-flux boundary condition. Secondly, the convergence at a belt of width $\sqrt{\epsilon}$ near $\partial\mathcal{M}$ is slower than the inner part of \mathcal{M} , see more in Lafon'06 paper.

3.4. Variance Term. Our purpose is to derive the large deviation bound for²

$$(164) \quad \text{Prob} \left(\frac{F}{G} - \frac{\mathbb{E}[F]}{\mathbb{E}[G]} \geq \alpha \right)$$

where $F = F(x_i) = \frac{1}{n} \sum_{j \neq i} k_\epsilon(x_i, x_j) f(x_j)$ and $G = G(x_i) = \frac{1}{n} \sum_{j \neq i} k_\epsilon(x, x_j)$. With x_1, x_2, \dots, x_n as i.i.d random variables, F and G are sample means (up to a scaling constant). Define a new random variable

$$Y = \mathbb{E}[G]F - \mathbb{E}[F]G - \alpha \mathbb{E}[G](G - \mathbb{E}[G])$$

which is of mean zero and Eqn. (164) can be rewritten as

$$\text{Prob}(Y \geq \alpha \mathbb{E}[G]^2).$$

²The opposite direction is omitted here.

For simplicity by *Markov (Chebyshev) inequality*³,

$$\text{Prob}(Y \geq \alpha \mathbb{E}[G]^2) \leq \frac{\mathbb{E}[Y^2]}{\alpha^2 \mathbb{E}[G]^4}$$

and setting the right hand side to be $\delta \in (0, 1)$, then with probability at least $1 - \delta$ the following holds

$$\alpha \leq \frac{\sqrt{\mathbb{E}[Y^2]}}{\mathbb{E}[G]^2 \sqrt{\delta}} \sim O\left(\frac{\sqrt{\mathbb{E}[Y^2]}}{\mathbb{E}[G]^2}\right).$$

It remains to bound

$$\begin{aligned} \mathbb{E}[Y^2] &= (\mathbb{E}G)^2 \mathbb{E}(F^2) - 2(\mathbb{E}G)(\mathbb{E}F)\mathbb{E}(FG) + (\mathbb{E}F)^2 \mathbb{E}(G^2) + \dots \\ &\quad + 2\alpha(\mathbb{E}G)[(\mathbb{E}F)\mathbb{E}(G^2) - (\mathbb{E}G)\mathbb{E}(FG)] + \alpha^2(\mathbb{E}G)^2(\mathbb{E}(G^2) - (\mathbb{E}G)^2). \end{aligned}$$

So it suffices to give $\mathbb{E}(F)$, $\mathbb{E}(G)$, $\mathbb{E}(FG)$, $\mathbb{E}(F^2)$, and $\mathbb{E}(G^2)$. The former two are given in bias and for the variance parts in latter three, let's take one simple example with $\mathbb{E}(G^2)$.

Recall that x_1, x_2, \dots, x_n are distributed i.i.d according to density $p(x)$, and

$$G(x) = \frac{1}{n} \sum_{j \neq i} k_\epsilon(x, x_j),$$

so

$$\text{Var}(G) = \frac{1}{n^2}(n-1) \left[\int_{\mathcal{M}} k_\epsilon(x, y))^2 p(y) dy - (\mathbb{E}k_\epsilon(x, y))^2 \right].$$

Look at the simplest case of 1-dimension flat \mathcal{M} for an illustrative example:

$$\int_{\mathcal{M}} (k_\epsilon(x, y))^2 p(y) dy = \int_{\mathbb{R}} \frac{1}{\sqrt{\epsilon}} h^2(z^2) (p(x) + p'(x)(\sqrt{\epsilon}z + O(\epsilon))) dz,$$

$$\text{let } M_2 = \int_{\mathbb{R}} h^2(z^2) dz$$

$$\int_{\mathcal{M}} (k_\epsilon(x, y))^2 p(y) dy = p(x) \cdot \frac{1}{\sqrt{\epsilon}} M_2 + O(\sqrt{\epsilon}).$$

Recall that $\mathbb{E}k_\epsilon(x, y) = O(1)$, we finally have

$$\text{Var}(G) \sim \frac{1}{n} \left[\frac{p(x)M_2}{\sqrt{\epsilon}} + O(1) \right] \sim \frac{1}{n\sqrt{\epsilon}}.$$

Generally, for d -dimensional case, $\text{Var}(G) \sim n^{-1}\epsilon^{-\frac{d}{2}}$. Similarly one can derive estimates on $\text{Var}(F)$.

Ignoring the joint effect of $\mathbb{E}(FG)$, one can somehow get a rough estimate based on $F/G = [\mathbb{E}(F) + O(\sqrt{\mathbb{E}(F^2)})]/[\mathbb{E}(G) + O(\sqrt{\mathbb{E}(G^2)})]$ where we applied the Markov inequality on both the numerator and denominator. Combining those estimates together, we have the following,

$$\begin{aligned} \frac{F}{G} &= \frac{fp + \epsilon \frac{m_2}{2}(\Delta(fp) + \mathbb{E}[fp]) + O(\epsilon^2, n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}})}{p + \epsilon \frac{m_2}{2}(\Delta p + \mathbb{E}[p]) + O(\epsilon^2, n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}})} \\ &= f + \epsilon \frac{m_2}{2}(\Delta p + \mathbb{E}[p]) + O(\epsilon^2, n^{-\frac{1}{2}}\epsilon^{-\frac{d}{4}}), \end{aligned}$$

³It means that $\text{Prob}(X > \alpha) \leq \mathbb{E}(X^2)/\alpha^2$. A Chernoff bound with exponential tail can be found in Singer'06.

here $O(B_1, B_2)$ denotes the dominating one of the two bounds B_1 and B_2 in the asymptotic limit. As a result, the error (bias + variance) of $L_{\epsilon,\alpha}$ (dividing another ϵ) is of the order

$$(165) \quad O(\epsilon, n^{-\frac{1}{2}} \epsilon^{-\frac{d}{4}-1}).$$

In [Sin06] paper, the last term in the last line is improved to

$$(166) \quad O(\epsilon, n^{-\frac{1}{2}} \epsilon^{-\frac{d}{4}-\frac{1}{2}}),$$

where the improvement is by carefully analyzing the large deviation bound of $\frac{F}{G}$ around $\frac{\mathbb{E} F}{\mathbb{E} G}$ shown above, making use of the fact that F and G are correlated. Technical details are not discussed here.

In conclusion, we need to choose ϵ to balance bias error and variance error to be both small. For example, by setting the two bounds in Eqn. (166) to be of the same order we have

$$\epsilon \sim n^{-1/2} \epsilon^{-1/2-d/4},$$

that is

$$\epsilon \sim n^{-1/(3+d/2)},$$

so the total error is $O(n^{-1/(3+d/2)})$.

4. *Vector Diffusion Map

In this class, we introduce the topic of vector Laplacian on graphs and vector diffusion map.

The ideas for vector Laplacian on graphs and vector diffusion mapping are a natural extension from graph Laplacian operator and diffusion mapping on graphs. The reason why diffusion mapping is important is that previous dimension reduction techniques, such as the PCA and MDS, ignore the intrinsic structure of the manifold. By contrast, diffusion mapping derived from graph Laplacian is the optimal embedding that preserves locality in a certain way. Moreover, diffusion mapping gives rise to a kind of metric called diffusion distance. Manifold learning problems involving vector bundle on graphs provide the demand for vector diffusion mapping. And since vector diffusion mapping is an extension from diffusion mapping, their properties and convergence behavior are similar.

The application of vector diffusion mapping is not restricted to manifold learning however. Due to its usage of optimal registration transformation, it is also a valuable tool for problems in computer vision and computer graphics, for example, optimal matching of 3D shapes.

The organization of this lecture notes is as follows: We first review graph Laplacian and diffusion mapping on graphs as the basis for vector diffusion mapping. We then introduce three examples of vector bundles on graphs. After that, we come to vector diffusion mapping. Finally, we introduce some conclusions about the convergence of vector diffusion mapping.

4.1. graph Laplacian and diffusion mapping.

4.2. graph Laplacian. The goal of graph Laplacian is to discover the intrinsic manifold structure given a set of data points in space. There are three steps of constructing the graph Laplacian operator:

- construct the graph using either the ϵ -neighborhood way (for any data point, connect it with all the points in its ϵ -neighborhood) or the k -nearest neighbor way (connect it with its k -nearest neighbors);
- construct the weight matrix. Here we can use the simple-minded binary weight (0 or 1), or use the heat kernel weight. For undirected graph, the weight matrix is symmetric;
- denote \mathcal{D} as the diagonal matrix with $\mathcal{D}(i, i) = \deg(i)$, $\deg(i) := \sum_j w_{ij}$. The graph Laplacian operator is:

$$L = \mathcal{D} - W$$

The graph Laplacian has the following properties:

- $\forall f : V \rightarrow \mathbb{R}$, $f^T L f = \sum_{(i,j) \in E} w_{ij} (f_i - f_j)^2 \geq 0$
- G is connected $\Leftrightarrow f^T L f > 0, \forall f \neq 0$, where $\vec{1} = (1, \dots, 1)^T$
- G has k -connected components $\Leftrightarrow \dim(\ker(L)) = k$
(this property is compatible with the previous one, since $L\vec{1} = 0$)
- Kirchhoff's Matrix Tree theorem:
Consider a connected graph G and the binary weight matrix: $w_{ij} = \begin{cases} 1, & (i, j) \in E \\ 0, & \text{otherwise} \end{cases}$, denote the eigenvalues of L as $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$, then $\#\{\text{T: T is a spanning tree of } G\} = \frac{1}{n} \lambda_2 \dots \lambda_n$
- Fieldler Theory, which will be introduced in later chapters.

We can have a further understanding of Graph Laplacian using the language of exterior calculus on graph.

We give the following denotations:

$V = \{1, 2, \dots, |V|\}$. \vec{E} is the oriented edge set that for $(i, j) \in E$ and $i < j$, $\langle i, j \rangle$ is the positive orientation, and $\langle j, i \rangle$ is the negative orientation.

$\delta_0 : \mathbb{R}^V \rightarrow \mathbb{R}^{\vec{E}}$ is a coboundary map, such that

$$\delta_0 \circ f(i, j) = \begin{cases} f_i - f_j, & \langle i, j \rangle \in \vec{E} \\ 0, & \text{otherwise} \end{cases}$$

It is easy to see that $\delta_0 \circ f(i, j) = -\delta_0 \circ f(j, i)$

The inner product of operators on $\mathbb{R}^{\vec{E}}$ is defined as:

$$\langle u, v \rangle = \sum_{i,j} w_{ij} u_{ij} v_{ij}$$

$$u^* := u \text{ diag}(w_{ij})$$

where $\text{diag}(w_{ij}) \in \mathbb{R}^{\frac{n(n-1)}{2} \times \frac{n(n-1)}{2}}$ is the diagonal matrix that has w_{ij} on the diagonal position corresponding to $\langle i, j \rangle$.

$$u^* v = \langle u, v \rangle$$

Then,

$$L = D - W = \delta_0^T \text{diag}(w_{ij}) \delta_0 = \delta_0^* \delta_0$$

We first look at the graph Laplacian operator. We solve the generalized eigenvalue problem:

$$Lf = \lambda \mathcal{D}f$$

denote the generalized eigenvalues as:

$$0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

and the corresponding generalized eigenvectors:

$$f_1, \dots, f_n$$

we have already obtained the m-dimensional Laplacian eigenmap:

$$\mathbf{x}_i \rightarrow (\mathbf{f}_1(\mathbf{i}), \dots, \mathbf{f}_m(\mathbf{i}))$$

We now explains that this is the optimal embedding that preserves locality in the sense that connected points stays as close as possible. Specifically speaking, for the one-dimensional embedding, the problem is:

$$\min \sum_{i,j} (y_i - y_j)^2 w_{ij} = 2 \min_{\mathbf{y}} \mathbf{y}^T L \mathbf{y}$$

$$y^T L y = y^T \mathcal{D}^{-\frac{1}{2}} (I - \mathcal{D}^{-\frac{1}{2}} W \mathcal{D}^{-\frac{1}{2}}) \mathcal{D}^{-\frac{1}{2}} y$$

Since $I - \mathcal{D}^{-\frac{1}{2}} W \mathcal{D}^{-\frac{1}{2}}$ is symmetric, the object is minimized when $\mathcal{D}^{-\frac{1}{2}} y$ is the eigenvector for the second smallest eigenvalue(the first smallest eigenvalue is 0) of $I - \mathcal{D}^{-\frac{1}{2}} W \mathcal{D}^{-\frac{1}{2}}$, which is the same with λ_2 , the second smallest generalized eigenvalue of L .

Similarly, the m-dimensional optimal embedding is given by $Y = (\mathbf{f}_1, \dots, \mathbf{f}_m)$.

In diffusion map, the weights are used to define a discrete random walk. The transition probability in a single step from i to j is:

$$a_{ij} = \frac{w_{ij}}{\deg(i)}$$

Then the transition matrix $A = \mathcal{D}^{-1} W$.

$$A = \mathcal{D}^{-\frac{1}{2}} (\mathcal{D}^{-\frac{1}{2}} W \mathcal{D}^{-\frac{1}{2}}) \mathcal{D}^{\frac{1}{2}}$$

Therefore, A is similar to a symmetric matrix, and has n real eigenvalues μ_1, \dots, μ_n and the corresponding eigenvectors ϕ_1, \dots, ϕ_n .

$$A \phi_i = \mu_i \phi_i$$

A^t is the transition matrix after t steps. Thus, we have:

$$A^t \phi_i = \mu_i^t \phi_i$$

Define Λ as the diagonal matrix with $\Lambda(i, i) = \mu_i$, $\Phi = [\phi_1, \dots, \phi_n]$. The diffusion map is given by:

$$\Phi_t := \Phi \Lambda^t = [\mu_1^t \phi_1, \dots, \mu_n^t \phi_n]$$

4.3. the embedding given by diffusion map. $\Phi_t(i)$ denotes the i th row of Φ_t .

$$\langle \Phi_t(i), \Phi(j) \rangle = \sum_{k=1}^n \frac{A^t(i, k)}{\sqrt{\deg(k)}} \frac{A^t(j, k)}{\sqrt{\deg(k)}}$$

we can thus define a distance called *diffusion distance*

$$d_{DM,t}^2(i, j) := \langle \Phi_t(i), \Phi(i) \rangle + \langle \Phi_t(j), \Phi(j) \rangle - 2\langle \Phi_t(i), \Phi(j) \rangle = \sum_{k=1}^n \frac{(A^t(i, k) - A^t(j, k))^2}{\deg(k)}$$

4.4. Examples of vector bundles on graph.

- (1) Wind velocity field on globe:

To simplify the problem, we consider the two dimensional mesh on the globe(the latitude and the longitude). Each node on the mesh has a vector \vec{f} which is the wind velocity at that place.

- (2) Local linear regression:

The goal of local linear regression is to give an approximation of the regression function at an arbitrary point in the variable space.

Given the data $(y_i, \vec{x}_i)_{i=1}^n$ and an arbitrary point $\vec{x}, \vec{x}, \vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^p$, we want to find $\vec{\beta} := (\beta_0, \beta_1, \dots, \beta_p)^T$ that minimize $\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 K_n(\vec{x}_i, \vec{x})$. Here $K_n(\vec{x}_i, \vec{x})$ is a kernel function that defines the weight for \vec{x}_i at the point \vec{x} . For example, we can use the Nadaraya-Watson kernel $K_n(\vec{x}_i, \vec{x}) = e^{-\frac{\|\vec{x} - \vec{x}_i\|^2}{n}}$.

For a graph $G=(V,E)$, each point $\vec{x} \in V$ has a corresponding vector $\vec{\beta}(\vec{x})$. We therefore get a vector bundle on the graph $G(V,E)$.

Here $\vec{\beta}$ is kind of a gradient. In fact, if y and \vec{x} has the relationship $y = f(\vec{x})$, then $\beta = (f(\vec{x}), \nabla f(\vec{x}))^T$.

- (3) Social networks:

If we see users as vertices and the relationship bonds that connected users as edges, then a social network naturally gives rise to a graph $G=(V,E)$. Each user has an attribute profile containing all kinds of personal information, and a certain kind of information can be described by a vector \vec{f} recording different aspects. Again, we get a vector bundle on graph.

4.5. optimal registration transformation. Like in graph eigenmap, we expect the embedding \vec{f} to be preserve locality to a certain extent, which means that we expect the embedding of connected points to be sufficiently close. In the graph Laplacian case, we use $\sum_{i \sim j} w_{ij} \|\vec{f}_i - \vec{f}_j\|^2$. However, for vector bundle on graphs, subtraction of vectors at different points may not be done directly due to the curvature of the manifold. What makes sense should be the difference of vectors compared with the tangent spaces at the certain points. Therefore, we borrow the idea of parallel transport from differential geometry. Denote O_{ij} as the parallel transport operator from the tangent space at x_j to the tangent space at x_i . We want to find out the embedding that minimizes

$$\sum_{i \sim j} w_{ij} \|\vec{f}_i - O_{ij} \vec{f}_j\|^2$$

we will later define the vector diffusion mapping, and using the similar argument as in diffusion mapping, it is easy to see that vector diffusion mapping gives the optimal embedding that preserves locality in this sense.

we now discuss how we get the approximation of parallel transport operator given the data set.

The approximation of the tangent space at a certain point x_i is given by local PCA. Choose ϵ_i to be sufficiently small, and denote $x_{i_1}, \dots, x_{i_{N_i}}$ as the data points in the ϵ_i -neighborhood of x_i . Define

$$X_i := [x_{i_1} - x_i, \dots, x_{i_{N_i}} - x_i]$$

Denote D_i as the diagonal matrix with

$$D_i(j, j) = \sqrt{K\left(\frac{\|x_{i_j} - x_i\|}{\epsilon_i}\right)}, \quad j = 1, \dots, N_i$$

$$B_i := X_i D_i$$

Perform SVD on B_i :

$$B_i = U_i \Sigma_i V_i^T$$

We use the first d columns of U_i (which are the left eigenvectors of the d largest eigenvalues of B_i) to form an approximation of the tangent space at x_i . That is,

$$O_i = [u_{i_1}, \dots, u_{i_d}]$$

Then O_i is a numerical approximation to an orthonormal basis of the tangent space at x_i .

For connected points x_i and x_j , since they are sufficiently close to each other, their tangent space should be close. Therefore, $O_i O_{ij}$ and O_j should also be close. We there use the closest orthogonal matrix to $O_i^T O_j$ as the approximation of the parallel transport operator from x_j to x_i :

$$\rho_{ij} := \operatorname{argmin}_{O \text{ orthogonal}} \|O - O_i^T O_j\|_{HS}$$

where $\|A\|_{HS}^2 = \operatorname{Tr}(AA^T)$ is the Hilbert-Schmidt norm.

4.6. Vector Laplacian. Given the weight matrix $W = (w_{ij})$, we denote

$$D := \begin{pmatrix} \deg(1)I_p & & \\ & \ddots & \\ & & \deg(n)I_p \end{pmatrix} \in \mathbb{R}^{np \times np}$$

where $\deg(i) = \sum_j w_{ij}$ as in graph Laplacian.

Define S as the block matrix with

$$S_{ij} = \begin{cases} w_{ij} \rho_{ij}, & i \sim j \\ 0, & \text{otherwise} \end{cases}$$

The vector Laplacian is then defined as $\mathcal{L} = D - S$

Like Graph Laplacian, we introduce an orientation on E and a coboundary map $\delta_0 : (\mathbb{R}^d)^V \rightarrow (\mathbb{R}^d)^{\vec{E}}$

$$\delta_0 \circ f(i, j) = \begin{cases} \vec{f}_i - \rho_{ij} \vec{f}_j, & \langle i, j \rangle \in \vec{E} \\ 0, & \text{otherwise} \end{cases}, \text{ where } f = (\vec{f}_1, \dots, \vec{f}_n)^T$$

Inner product on $(\mathbb{R}^d)^{\vec{E}}$ is defined as

$$\langle u, v \rangle = \sum_{i,j} w_{ij} u_{ij}^T v_{ij}$$

$$u^* := u \operatorname{diag}(w_{ij}), u^* v = \langle u, v \rangle$$

If we let ρ_{ij} be orthogonal, $\forall i, j, s.t. \langle i, j \rangle \in \vec{E}$, then, $L = D - W = \delta_0^T \operatorname{diag}(w_{ij}) \delta_0 = \delta_0^* \delta_0$.

Analogous properties with Graph Laplacian:

- G has k connected components $\Leftrightarrow \dim \ker(\mathcal{L}) = kp$
- generalized Matrix tree theorem.

4.7. Vector diffusion mapping.

$$\begin{aligned} \mathcal{L} &= D - S = D(I - D^{-1}S) \\ D^{-1}S &= D^{-\frac{1}{2}}SD^{-\frac{1}{2}} \end{aligned}$$

Denote

$$\tilde{S} := D^{-\frac{1}{2}}SD^{-\frac{1}{2}}$$

\tilde{S} has nd real eigenvalues $\lambda_1, \dots, \lambda_{nd}$ and the corresponding eigenvectors v_1, \dots, v_{nd} . Thinking of these vectors of length nd in blocks of d, we denote $v_k(i)$ as the ith block of v_k .

The spectral decompositions of $\tilde{S}(i, j)$ and $\tilde{S}^{2t}(i, j)$ are given by:

$$\begin{aligned} \tilde{S}(i, j) &= \sum_{k=1}^{nd} \lambda_k v_k(i) v_k(j)^T \\ \therefore \tilde{S}^{2t}(i, j) &= \sum_{k=1}^{nd} \lambda_k^{2t} v_k(i) v_k(j)^T \end{aligned}$$

We use $\|\tilde{S}^{2t}(i, j)\|_{HS}^2$ to measure the affinity between i and j. Thus,

$$\begin{aligned} \|\tilde{S}^{2t}(i, j)\|_{HS}^2 &= \operatorname{Tr}(\tilde{S}^{2t}(i, j) \tilde{S}^{2t}(i, j)^T) \\ &= \sum_{k,l=1}^{nd} (\lambda_k \lambda_l)^{2t} \operatorname{Tr}(v_k(i) v_k(j)^T v_l(j) v_l(i)^T) \\ &= \sum_{k,l=1}^{nd} (\lambda_k \lambda_l)^{2t} \operatorname{Tr}(v_k(j)^T v_l(j) v_l(i)^T v_k(i)) \\ &= \sum_{k,l=1}^{nd} (\lambda_k \lambda_l)^{2t} \langle v_k(j), v_l(j) \rangle \langle v_k(i), v_l(i) \rangle \end{aligned}$$

The vector diffusion mapping is defined as:

$$V_t : i \rightarrow ((\lambda_k \lambda_l)^t \langle v_k(i), v_l(i) \rangle)_{k,l=1}^{nd}$$

Like graph Laplacian, $\|\tilde{S}^{2t}(i, j)\|_{HS}^2$ is actually an inner product:

$$\|\tilde{S}^{2t}(i, j)\|_{HS}^2 = \langle V_t(i), V_t(j) \rangle$$

This gives rise to a distance called vector diffusion distance:

$$d_{VDM,t}^2 = \langle V_t(i), V_t(i) \rangle + \langle V_t(j), V_t(j) \rangle - 2 \langle V_t(i), V_t(j) \rangle$$

4.8. Normalized Vector Diffusion Mappings. An important kind of normalized VDM is obtained as follows:

Take $0 \leq \alpha \leq 1$,

$$\begin{aligned} W_\alpha &:= \mathcal{D}^{-\alpha} W \mathcal{D}^{-\alpha} \\ S_\alpha &:= D^{-\alpha} S D^{-\alpha} \\ \deg_\alpha(i) &:= \sum_{j=1}^n W_\alpha(i, j) \end{aligned}$$

We define $\mathcal{D}_\alpha \in \mathbb{R}^{n \times n}$ as the diagonal matrix with

$$\mathcal{D}_\alpha(i, i) = \deg_\alpha(i)$$

and $D_\alpha \in \mathbb{R}^{nd \times nd}$ as the block diagonal matrix with

$$D_\alpha(i, i) = \deg_\alpha(i) I_d$$

We can then get the vector diffusion mapping $V_{\alpha,t}$ using S_α and D_α instead of S and D.

4.9. Convergence of VDM. We first introduce some concepts.

Suppose \mathcal{M} is a smooth manifold, and $T_{\mathcal{M}}$ is a *tensor bundle* on \mathcal{M} . When the rank of $T_{\mathcal{M}}$ is 0, it is the set of functions on \mathcal{M} . When the rank of $T_{\mathcal{M}}$ is 1, it is the set of vector fields on \mathcal{M} .

The *connection Laplacian operator* is:

$$\nabla_{X,Y}^2 T = -(\nabla_X \nabla_Y T - \nabla_{\nabla_X Y} T)$$

where $\nabla_X Y$ is the covariant derivative of Y over X.

Intuitively, we can see the first item of the connection Laplacian operator as the sum of the change of T over X and over Y, and the second item as the overlapped part of the change of T over X and over Y. The remainder can be seen as an operator that differentiates the vector fields in the direction of two orthogonal vector fields.

Now we introduce some results about convergence.

The normalized graph Laplacian converges to the Laplace-Beltrami operator:

$$(\mathcal{D}^{-1} W - I)f \rightarrow c\Delta f$$

for sufficiently smooth f and some constant c.

For VDM, $D_\alpha^{-1} S_\alpha - I$ converges to the connection Laplacian operator [SW12] plus some potential terms. When $\alpha = 1$, $D_1^{-1} S_1 - I$ converges to exactly the connection Laplacian operator:

$$(D_1^{-1} S_1 - I)X \rightarrow c\nabla^2 X$$

CHAPTER 8

Semi-supervised Learning

1. Introduction

Problem: $x_1, x_2, \dots, x_l \in V_l$ are labeled data, that is data with the value $f(x_i), f \in V \rightarrow \mathbb{R}$ observed. $x_{l+1}, x_{l+2}, \dots, x_{l+u} \in V_u$ are unlabeled. Our concern is how to fully exploiting the information (like geometric structure in distribution) provided in the labeled and unlabeled data to find the unobserved labels.

This kind of problem may occur in many situations, like ZIP Code recognition. We may only have a part of digits labeled and our task is to label the unlabeled ones.

2. Harmonic Extension of Functions on Graph

Suppose the whole graph is $G = (V, E, W)$, where $V = V_l \cup V_u$ and weight matrix is partitioned into blocks $W = \begin{pmatrix} W_{ll} & W_{lu} \\ W_{ul} & W_{uu} \end{pmatrix}$. As before, we define $D = \text{diag}(d_1, d_2, \dots, d_n) = \text{diag}(D_l, D_u)$, $d_i = \sum_{j=1}^n W_{ij}$, $L = D - W$. The goal is to find $f_u = (f_{l+1}, \dots, f_{l+u})^T$ such that

$$\begin{aligned} \min \quad & f^T L f \\ \text{s.t.} \quad & f(V_l) = f_l \end{aligned}$$

where $f = \begin{pmatrix} f_l \\ f_u \end{pmatrix}$. Note that

$$f^T L f = (f_l^T, f_u^T) L \begin{pmatrix} f_l \\ f_u \end{pmatrix} = f_u^T L_{uu} f_u + f_l^T L_{ll} f_l + 2f_u^T L_{ul} f_l$$

So we have:

$$\frac{\partial f^T L f}{\partial f_u} = 0 \Rightarrow 2L_{uu}f_u + 2L_{lu}f_u = 0 \Rightarrow f_u = -L_{uu}^{-1}L_{ul}f_l = (D_u - W_{uu})^{-1}W_{ul}f_l$$

3. Explanation from Gaussian Markov Random Field

If we consider $f : V \rightarrow \mathbb{R}$ are Gaussian random variables on graph nodes whose inverse covariance matrix (precision matrix) is given by unnormalized graph Laplacian L (sparse but singular), i.e. $f \sim \mathcal{N}(0, \Sigma)$ where $\Sigma^{-1} = L$ (interpreted as a pseudo inverse). Then the conditional expectation of f_u given f_l is:

$$f_u = \Sigma_{ul}\Sigma_{ll}^{-1}f_l$$

where

$$\Sigma = \begin{bmatrix} \Sigma_{ll} & \Sigma_{lu} \\ \Sigma_{ul} & \Sigma_{uu} \end{bmatrix}$$

Block matrix inversion formula tells us that when A and D are invertible,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \cdot \begin{bmatrix} X & Y \\ Z & W \end{bmatrix} = I \Rightarrow \begin{bmatrix} X & Y \\ Z & W \end{bmatrix} = \begin{bmatrix} S_D^{-1} & -A^{-1}BS_A^{-1} \\ -D^{-1}CS_D^{-1} & S_A^{-1} \end{bmatrix}$$

$$\begin{bmatrix} X & Y \\ Z & W \end{bmatrix} \cdot \begin{bmatrix} A & B \\ C & D \end{bmatrix} = I \Rightarrow \begin{bmatrix} X & Y \\ Z & W \end{bmatrix} = \begin{bmatrix} S_D^{-1} & -S_D^{-1}BD^{-1} \\ -S_A^{-1}CA^{-1} & S_A^{-1} \end{bmatrix}$$

where $S_A = D - CA^{-1}B$ and $S_D = A - BD^{-1}C$ are called Schur complements of A and D , respectively. The matrix expressions for inverse are equivalent when the matrix is invertible.

The graph Laplacian

$$L = \begin{bmatrix} D_l - W_{ll} & -W_{lu} \\ -W_{ul} & D_u - W_{uu} \end{bmatrix}$$

is not invertible. $D_l - W_{ll}$ and $D_u - W_{uu}$ are both strictly diagonally dominant, i.e. $D_l(i, i) > \sum_j |W_{ll}(i, j)|$, whence they are invertible by Gershgorin Circle Theorem. However their Schur complements $S_{D_u - W_{uu}}$ and $S_{D_l - W_{ll}}$ are still not invertible and the block matrix inversion formula above can not be applied directly. To avoid this issue, we define a regularized version of graph Laplacian

$$L_\lambda = L + \lambda I, \quad \lambda > 0$$

and study its inverse $\Sigma_\lambda = L_\lambda^{-1}$.

By the block matrix inversion formula, we can set Σ as its right inverse above,

$$\Sigma_\lambda = \begin{bmatrix} S_{\lambda+D_u-W_{uu}}^{-1} & (\lambda + D_l - W_{ll})^{-1}W_{lu}S_{\lambda+D_l-W_{ll}}^{-1} \\ (\lambda + D_u - W_{uu})^{-1}W_{ul}S_{\lambda+D_u-W_{uu}}^{-1} & S_{\lambda+D_l-W_{ll}}^{-1} \end{bmatrix}$$

Therefore,

$$f_{u,\lambda} = \Sigma_{ul,\lambda}\Sigma_{ll,\lambda}^{-1}f_l = (\lambda + D_u - W_{uu})^{-1}W_{ul}f_l,$$

whose limit however exists $\lim_{\lambda \rightarrow 0} f_{u,\lambda} = (D_u - W_{uu})^{-1}W_{ul}f_l = f_u$. This implies that f_u can be regarded as the conditional mean given f_l .

4. Explanation from Transition Path Theory

We can also view the problem as a random walk on graph. Constructing a graph model with transition matrix $P = D^{-1}W = \begin{pmatrix} P_{ll} & P_{lu} \\ P_{ul} & P_{uu} \end{pmatrix}$. Assume that the labeled data are binary (classification). That is, for $x_i \in V_l$, $f(x_i) = 0$ or 1. Denote

- $V_0 = \{i \in V_l : f_i = f(x_i) = 0\}$
- $V_1 = \{i \in V_l : f_i = f(x_i) = 1\}$
- $V = V_0 \cup V_1 \cup V_u$ where $V_l = V_0 \cup V_1$

With this random walk on graph P , f_u can be interpreted as hitting time or first passage time of V_1 .

PROPOSITION 4.1. Define hitting time

$$\tau_i^k = \inf\{t \geq 0 : x(0) = i, x(t) \in V_k\}, \quad k = 0, 1$$

Then for $\forall i \in V_u$,

$$f_i = \text{Prob}(\tau_i^1 < \tau_i^0)$$

i.e.

$$f_i = \text{Prob}(\text{trajectory starting from } x_i \text{ hit } V_1 \text{ before } V_0)$$

Note that the probability above also called committor function in Transition Path Theory of Markov Chains.

PROOF. Define the committor function,

$$q_i^+ = \text{Prob}(\tau_i^1 < \tau_i^0) = \begin{cases} 1 & x_i \in V_1 \\ 0 & x_i \in V_0 \\ \sum_{j \in V} P_{ij} q_j^+ & i \in V_u \end{cases}$$

This is because $\forall i \in V_u$,

$$\begin{aligned} q_i^+ &= \Pr(\tau_{iV_1} < \tau_{iV_0}) \\ &= \sum_j P_{ij} q_j^+ \\ &= \sum_{j \in V_1} P_{ij} q_j^+ + \sum_{j \in V_0} P_{ij} q_j^+ + \sum_{j \in V_u} P_{ij} q_j^+ \\ &= \sum_{j \in V_1} P_{ij} + \sum_{j \in V_u} P_{ij} q_j^+ \end{aligned}$$

$$\therefore q_u^+ = P_{ul} f_l + P_{uu} q_u^+ = D_u^{-1} W_{ul} f_l + D_u^{-1} W_{uu} q_u^+$$

multiply D_u to both side and reorganize:

$$(D_u - W_{uu}) q_u^+ = W_{ul} f_l$$

If $D_u - W_{uu}$ is reversible, we get:

$$q_u^+ = (D_u - W_{uu})^{-1} W_{ul} f_l = f_u$$

i.e. f_u is the committor function on V_u . □

The result coincides with we obtained through the view of gaussian markov random field.

5. Well-posedness

One natural problem is: if we only have a fixed amount of labeled data, can we recover labels of an infinite amount of unobserved data? This is called well-posedness. [Nadler-Srebro 2009] gives the following result:

- If $x_i \in \mathbb{R}^1$, the problem is well-posed.
- If $x_i \in \mathbb{R}^d (d \geq 3)$, the problem is ill-posed in which case $D_u - W_{uu}$ becomes singular and f becomes a bump function (f_u is almost always zeros or ones except on some singular points).

Here we can give a brief explanation:

$$f^T L f \sim \int \|\nabla f\|_2$$

If we have $V_l = \{0, 1\}$, $f(x_0) = 0$, $f(x_1) = 1$ and let $f_\epsilon(x) = \begin{cases} \frac{\|x-x_0\|_2^2}{\epsilon^2} & \|x-x_0\|_2 < \epsilon \\ 1 & \text{otherwise} \end{cases}$.

From multivariable calculus,

$$\int \|\nabla f\|_2 = c\epsilon^{d-2}.$$

Since $d \geq 3$, so $\epsilon \rightarrow 0 \Rightarrow \int \|\nabla f\|_2 \rightarrow 0$. So $f_\epsilon(x)$ ($\epsilon \rightarrow 0$) converges to a bump function which is one almost everywhere except x_0 whose value is 0. No generalization ability is learned for such bump functions.

This means in high dimensional case, to obtain a smooth generalization, we have to add constraints more than the norm of the first order derivatives. We also have a theorem to illustrate what kind of constraint is enough for a good generalization:

THEOREM 5.1 (Sobolev embedding Theorem). $f \in \mathbf{W}^{s,p}(\mathbb{R}^d) \iff f$ has s 'th order weak derivative $f^{(s)} \in \mathbf{L}_p$,

$$s > \frac{d}{2} \Rightarrow \mathbf{W}^{s,2} \hookrightarrow \mathbf{C}(\mathbb{R}^d).$$

So in \mathbb{R}^d , to obtain a continuous function, one needs smoothness regularization $\int \|\nabla^s f\|$ with degree $s > d/2$. To implement this in discrete Laplacian setting, one may consider iterative Laplacian L^s which might converge to high order smoothness regularization.

CHAPTER 9

Beyond graphs: high dimensional topological/geometric analysis

1. From Graph to Simplicial Complex

DEFINITION (Simplicial Complex). An abstract simplicial complex is a collection Σ of subsets of V which is closed under inclusion (or deletion), i.e. $\tau \in \Sigma$ and $\sigma \subseteq \tau$, then $\sigma \in \Sigma$.

We have the following examples:

- Chess-board Complex
- Point cloud data:
 - Nerve complex
 - Cech, Rips, Witness complex
 - Mayer-Vietoris Blowup
- Term-document cooccurrence complex
- Clique complex in pairwise comparison graphs
- Strategic complex in flow games

EXAMPLE (Chess-board Complex). Let V be the positions on a Chess board. Σ collects position subsets of V where one can place queens (rooks) without capturing each other. It is easy to check the closedness under deletion: if $\sigma \in \Sigma$ is a set of “safe” positions, then any subset $\tau \subseteq \sigma$ is also a set of “safe” positions

EXAMPLE (Nerve Complex). Define a cover of X , $X = \cup_{\alpha} U_{\alpha}$. $V = \{U_{\alpha}\}$ and define $\Sigma = \{U_I : \cap_{\alpha \in I} U_{\alpha} \neq \emptyset\}$.

- Closedness under deletion
- Can be applied to any topological space X
- In a metric space (X, d) , if $U_{\alpha} = B_{\epsilon}(t_{\alpha}) := \{x \in X : d(x - t_{\alpha}) \leq \epsilon\}$, we have **Cech complex** C_{ϵ} .
- **Nerve Theorem**: if every U_I is contractible, then X has the same homotopy type as Σ .
- Cech complex is hard to compute, even in Euclidean space
- One can easily compute an upper bound for Cech complex

Construct a Cech subcomplex of 1-dimension, i.e. a graph with edges connecting point pairs whose distance is no more than ϵ .

Find the clique complex, i.e. maximal complex whose 1-skeleton is the graph above, where every k -clique is regarded as a $k - 1$ simplex

EXAMPLE (Vietoris-Rips Complex). Let $V = \{x_{\alpha} \in X\}$. Define $VR_{\epsilon} = \{U_I \subseteq V : d(x_{\alpha}, x_{\beta}) \leq \epsilon, \alpha, \beta \in I\}$.

- Rips is easier to compute than Cech

even so, Rips is exponential to dimension generally

- However Vietoris-Rips CAN NOT preserve the homotopy type as Cech
- But there is still a hope to find a **lower bound** on homology –

THEOREM 1.1 (“Sandwich”).

$$VR_\epsilon \subseteq C_\epsilon \subseteq VR_{2\epsilon}$$

- If a homology group “persists” through $R_\epsilon \rightarrow R_{2\epsilon}$, then it must exists in C_ϵ ; but not the vice versa.
- All above gives rise to a filtration of simplicial complex

$$\emptyset = \Sigma_0 \subseteq \Sigma_1 \subseteq \Sigma_2 \subseteq \dots$$

- Functoriality of inclusion: there are homomorphisms between homology groups

$$0 \rightarrow H_1 \rightarrow H_2 \rightarrow \dots$$

- A persistent homology is the image of H_i in H_j with $j > i$.

EXAMPLE (Strong Witness Complex). Let $V = \{t_\alpha \in X\}$. Define $W_\epsilon^s = \{U_I \subseteq V : \exists x \in X, \forall \alpha \in I, d(x, t_\alpha) \leq d(x, V) + \epsilon\}$.

EXAMPLE (Week Witness Complex). Let $V = \{t_\alpha \in X\}$. Define $W_\epsilon^w = \{U_I \subseteq V : \exists x \in X, \forall \alpha \in I, d(x, t_\alpha) \leq d(x, V) + \epsilon\}$.

- V can be a set of landmarks, much smaller than X
- Monotonicity: $W_\epsilon^* \subseteq W_{\epsilon'}^*$ if $\epsilon \leq \epsilon'$
- But not easy to control homotopy types between W^* and X

EXAMPLE (Term-Document Occurrence complex, [LK10]). Left is a term-document co-occurrence matrix; Right is a simplicial complex representation of terms. Connectivity analysis captures more information than Latent Semantic Index.

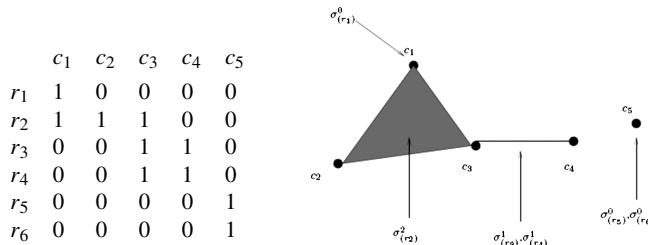


FIGURE 1. Term-Document Occurrence complex

EXAMPLE (Flag Complex of Paired Comparison Graph, Jiang-Lim-Yao-Ye 2011[JLY11]). Let V be a set of alternatives to be compared and undirected pair $(i, j) \in E$ if the pair is comparable. A flag complex χ_G consists all cliques as simplices or faces (e.g. 3-cliques as 2-faces and $k+1$ -cliques as k -faces), also called clique complex of G .

EXAMPLE (Strategic Simplicial Complex for Flow Games, Candogan-Menache-Ozdaglar-Parrilo 2011 [**CMOP11**]). Strategic simplicial complex is the clique complex of pairwise comparison graph $G = (V, E)$ of strategic profiles, where V consists of all strategy profiles of players and a pair of strategy $(x, x') \in E$ is comparable if only one player changes strategy from x to x' . Every finite game can be decomposed as the direct sum of potential games and zero-sum games (harmonic games).

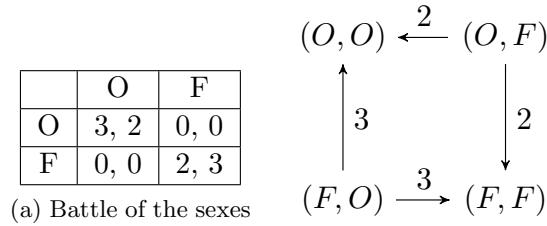


FIGURE 2. Illustration of Game Strategic Complex: Battle of Sex

2. Persistent Homology and Discrete Morse Theory

Recall that

THEOREM 2.1 (“Sandwich”).

$$VR_\epsilon \subseteq C_\epsilon \subseteq VR_{2\epsilon}$$

- If a homology group “persists” through $R_\epsilon \rightarrow R_{2\epsilon}$, then it must exist in C_ϵ ; but not the vice versa.
- All above gives rise to a filtration of simplicial complex

$$\emptyset = \Sigma_0 \subseteq \Sigma_1 \subseteq \Sigma_2 \subseteq \dots$$

- Functoriality of inclusion: there are homomorphisms between homology groups

$$0 \rightarrow H_1 \rightarrow H_2 \rightarrow \dots$$

- A persistent homology is the image of H_i in H_j with $j > i$.

Persistent Homology is firstly proposed by Edelsbrunner-Letscher-Zomorodian, with an algebraic formulation by Zomorodian-Carlsson. The algorithm is equivalent to Robin Forman’s discrete Morse theory.

to be continued...

3. Exterior Calculus on Complex and Combinatorial Hodge Theory

We are going to study functions on simplicial complex, $l^2(V^d)$.

A basis of “forms”:

- $l^2(V)$: e_i ($i \in V$), so $f \in l^2(V)$ has a representation $f = \sum_{i \in V} f_i e_i$, e.g. global ranking score on V .
- $l^2(V^2)$: $e_{ij} = -e_{ji}$, $f = \sum_{(i,j)} f_{ij} e_{ij}$ for $f \in l^2(V^2)$, e.g. paired comparison scores on V^2 .
- $l^2(V^3)$: $e_{ijk} = e_{jki} = e_{kij} = -e_{jik} = -e_{kji} = -e_{ikj}$, $f = \sum_{ijk} f_{ijk} e_{ijk}$

- $l^2(V^{d+1})$: e_{i_0, \dots, i_d} is an alternating d -form

$$e_{i_0, \dots, i_d} = \text{sign}(\sigma) e_{\sigma(i_0), \dots, \sigma(i_d)},$$

where $\sigma \in \mathfrak{S}_d$ is a permutation on $\{0, \dots, d\}$.

Vector spaces of functions $l^2(V^{d+1})$ represented on such basis with an inner product defined, are called d -forms (cochains).

EXAMPLE. In the crowdsourcing ranking of world universities,

<http://www.allourideas.org/worldcollege/>,

V consists of world universities, E are university pairs in comparison, $l^2(V)$ consists of ranking scores of universities, $l^2(V^2)$ is made up of paired comparison data.

Discrete differential operators: k -dimensional **coboundary maps** $\delta_k : L^2(V^k) \rightarrow L^2(V^{k+1})$ are defined as the **alternating difference** operator

$$(\delta_k u)(i_0, \dots, i_{k+1}) = \sum_{j=0}^{k+1} (-1)^{j+1} u(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_{k+1})$$

- δ_k plays the role of **differentiation**
- $\delta_{k+1} \circ \delta_k = 0$

So we have chain map

$$L^2(V) \xrightarrow{\delta_0} L^2(V^2) \xrightarrow{\delta_1} L^2(V^3) \rightarrow \dots L^2(V^k) \xrightarrow{\delta_{k-1}} L^2(V^{k+1}) \xrightarrow{\delta_k} \dots$$

with $\delta_k \circ \delta_{k-1} = 0$.

EXAMPLE (Gradient, Curl, and Divergence). We can define discrete gradient and curl, as well as their adjoints

- $(\delta_0 v)(i, j) = v_j - v_i =: (\text{grad } v)(i, j)$
- $(\delta_1 w)(i, j, k) = (\pm)(w_{ij} + w_{jk} + w_{ki}) =: (\text{curl } w)(i, j, k)$, which measures the total flow-sum along the loop $i \rightarrow j \rightarrow k \rightarrow i$ and $(\delta_1 w)(i, j, k) = 0$ implies the paired comparison data is **path-independent**, which defines the **triangular transitivity subspace**
- for each alternative $i \in V$, the **combinatorial divergence**

$$(\text{div } w)(i) := -(\delta_0^T w)(i) := \sum w_{i*}$$

which measures the **inflow-outflow sum** at i and $(\delta_0^T w)(i) = 0$ implies alternative i is preference-neutral in all pairwise comparisons as a cyclic ranking passing through alternatives.

DEFINITION (Combinatorial Hodge Laplacian). Define the k -dimensional **combinatorial Laplacian**, $\Delta_k : L^2(V^{k+1}) \rightarrow L^2(C^{k+1})$ by

$$\Delta_k = \delta_{k-1} \delta_{k-1}^T + \delta_k^T \delta_k, \quad k > 0$$

- $k = 0$, $\Delta_0 = \delta_0^T \delta_0$ is the well-known **graph Laplacian**
- $k = 1$,

$$\Delta_1 = \text{curl} \circ \text{curl}^* - \text{div} \circ \text{grad}$$

- Important Properties:

Δ_k positive semi-definite

$\ker(\Delta_k) = \ker(\delta_{k-1}^T) \cap \ker(\delta_k)$: **k -Harmonics**, dimension equals to k -th Betti number

Hodge Decomposition Theorem

THEOREM 3.1 (Hodge Decomposition). The space of k -forms (cochains) $C^k(\mathcal{K}(G), \mathbb{R})$, admits an orthogonal decomposition into three

$$C^k(\mathcal{K}(G), \mathbb{R}) = \text{im}(\delta_{k-1}) \oplus H_k \oplus \text{im}(\delta_k^T)$$

where

$$H_k = \ker(\delta_{k-1}) \cap \ker(\delta_k^T) = \ker(\Delta_k).$$

- $\dim(H_k) = \beta_k$.

A simple understanding is possible via **Dirac** operator:

$$D = \delta + \delta^* : \oplus_k L^2(V^k) \rightarrow \oplus_k L^2(V^k)$$

Hence $D = D^*$ is self-adjoint. Combine the chain map

$$L^2(V) \xrightarrow{\delta_0} L^2(V^2) \xrightarrow{\delta_1} L^2(V^3) \rightarrow \dots L^2(V^k) \xrightarrow{\delta_{k-1}} L^2(V^{k+1}) \xrightarrow{\delta_k} \dots$$

into a big operator: Dirac operator.

Abstract Hodge Laplacian:

$$\Delta = D^2 = \delta\delta^* + \delta^*\delta,$$

since $\delta^2 = 0$.

By the Fundamental Theorem of Linear Algebra (Closed Range Theorem in Banach Space),

$$\oplus_k L^2(V^k) = \text{im}(D) \oplus \ker(D)$$

where

$$\text{im}(D) = \text{im}(\delta) \oplus \text{im}(\delta^*)$$

and $\ker(D) = \ker(\Delta)$ is the space of harmonic forms.

4. Applications of Hodge Theory: Statistical Ranking

4.1. HodgeRank on Graphs. Let $\wedge = \{1, \dots, m\}$ be a set of participants and $V = \{1, \dots, n\}$ be the set of videos to be ranked. Paired comparison data is collected as a function on $\wedge \times V \times V$, which is *skew-symmetric* for each participant α , *i.e.*, $Y_{ij}^\alpha = -Y_{ji}^\alpha$ representing the degree that α prefers i to j . The simplest setting is the binary choice, where

$$Y_{ij}^\alpha = \begin{cases} 1 & \text{if } \alpha \text{ prefers } i \text{ to } j, \\ -1 & \text{otherwise.} \end{cases}$$

In general, Y_{ij}^α can be used to represent paired comparison grades, *e.g.*, $Y_{ij}^\alpha > 0$ refers to the degree that α prefers i to j and the vice versa $Y_{ji}^\alpha = -Y_{ij}^\alpha < 0$ measures the dispreference degree [JYY11].

In this paper we shall focus on the binary choice, which is the simplest setting and the data collected in this paper belongs to this case. However the theory can be applied to the more general case with multiple choices above.

Such paired comparison data can be represented by a directed graph, or hypergraph, with n nodes, where each directed edge between i and j refers the preference indicated by Y_{ij}^α .

A nonnegative weight function $\omega : \wedge \times V \times V \rightarrow [0, \infty)$ is defined as,

$$(167) \quad \omega_{ij}^\alpha = \begin{cases} 1 & \text{if } \alpha \text{ makes a comparison for } \{i, j\}, \\ 0 & \text{otherwise.} \end{cases}$$

It may reflect the confidence level that a participant compares $\{i, j\}$ by taking different values, and this is however not pursued in this paper.

Our statistical rank aggregation problem is to look for some global ranking score $s : V \rightarrow R$ such that

$$(168) \quad \min_{s \in R^{|V|}} \sum_{i,j,\alpha} \omega_{ij}^\alpha (s_i - s_j - Y_{ij}^\alpha)^2,$$

which is equivalent to the following weighted least square problem

$$(169) \quad \min_{s \in R^{|V|}} \sum_{i,j} \omega_{ij} (s_i - s_j - \hat{Y}_{ij})^2,$$

where $\hat{Y}_{ij} = (\sum_\alpha \omega_{ij}^\alpha Y_{ij}^\alpha) / (\sum_\alpha \omega_{ij}^\alpha)$ and $\omega_{ij} = \sum_\alpha \omega_{ij}^\alpha$. For the principles behind such a choice, readers may refer [JLYY11].

A graph structure arises naturally from ranking data as follows. Let $G = (V, E)$ be a paired ranking graph whose vertex set is V , the set of videos to be ranked, and whose edge set is E , the set of video pairs which receive some comparisons, *i.e.*,

$$(170) \quad E = \left\{ \{i, j\} \in \binom{V}{2} \mid \sum_\alpha \omega_{i,j}^\alpha > 0 \right\}.$$

A pairwise ranking is called *complete* if each participant α in \wedge gives a total judgment of all videos in V ; otherwise it is called *incomplete*. It is called *balanced* if the paired comparison graph is k -regular with equal weights $\omega_{ij} = \sum_\alpha \omega_{ij}^\alpha \equiv c$ for all $\{i, j\} \in E$; otherwise it is called *imbalanced*. A complete and balanced ranking induces a complete graph with equal weights on all edges. The existing paired comparison methods in VQA often assume complete and balanced data. However, this is an unrealistic assumption for real world data, *e.g.* randomized experiments. Moreover in crowdsourcing, raters and videos come in an unspecified way and it is hard to control the test process with precise experimental designs. Nevertheless, as to be shown below, it is efficient to utilize some random sampling design based on random graph theory where for each participant a fraction of video pairs are chosen randomly. The HodgeRank approach adopted in this paper enables us a unified scheme which can deal with incomplete and imbalanced data emerged from random sampling in paired comparisons.

The minimization problem (169) can be generalized to a family of *linear models* in paired comparison methods [Dav88]. To see this, we first rewrite (169) in another simpler form. Assume that for each edge as video pair $\{i, j\}$, the number of comparisons is n_{ij} , among which a_{ij} participants have a preference on i over j (a_{ji} carries the opposite meaning). So $a_{ij} + a_{ji} = n_{ij}$ if no tie occurs. Therefore, for each edge $\{i, j\} \in E$, we have a preference probability estimated from data $\hat{\pi}_{ij} = a_{ij} / n_{ij}$. With this definition, the problem (169) can be rewritten as

$$(171) \quad \min_{s \in R^{|V|}} \sum_{\{i,j\} \in E} n_{ij} (s_i - s_j - (2\hat{\pi}_{ij} - 1))^2,$$

since $\hat{Y}_{ij} = (a_{ij} - a_{ji}) / n_{ij} = 2\hat{\pi}_{ij} - 1$ due to Equation (167).

General *linear models*, which are firstly formulated by G. Noether [Noe60], assume that the true preference probability can be fully decided by a linear scaling function on V , *i.e.*,

$$(172) \quad \pi_{ij} = \text{Prob}\{i \text{ is preferred over } j\} = F(s_i^* - s_j^*),$$

for some $s^* \in R^{|V|}$. F can be chosen as any symmetric cumulated distributed function. When only an empirical preference probability $\hat{\pi}_{ij}$ is observed, we can

map it to a skew-symmetric function by the inverse of F ,

$$(173) \quad \hat{Y}_{ij} = F^{-1}(\hat{\pi}_{ij}),$$

where $\hat{Y}_{ij} = -\hat{Y}_{ji}$. However, in this case, one can only expect that

$$(174) \quad \hat{Y}_{ij} = s_i^* - s_j^* + \varepsilon_{ij},$$

where ε_{ij} accounts for the noise. The case in (171) takes a linear F and is often called a *uniform model*. Below we summarize some well known models which have been studied extensively in [Dav88].

1. *Uniform* model:

$$(175) \quad \hat{Y}_{ij} = 2\hat{\pi}_{ij} - 1.$$

2. *Bradley-Terry* model:

$$(176) \quad \hat{Y}_{ij} = \log \frac{\hat{\pi}_{ij}}{1 - \hat{\pi}_{ij}}.$$

3. *Thurstone-Mosteller* model:

$$(177) \quad \hat{Y}_{ij} = F^{-1}(\hat{\pi}_{ij}).$$

where F is essentially the Gauss error function

$$(178) \quad F(x) = \frac{1}{\sqrt{2\pi}} \int_{-x/[2\sigma^2(1-\rho)]^{1/2}}^{\infty} e^{-\frac{1}{2}t^2} dt.$$

Note that constants σ and ρ will only contribute to a rescaling of the solution of (169).

4. *Angular transform* model:

$$(179) \quad \hat{Y}_{ij} = \arcsin(2\hat{\pi}_{ij} - 1).$$

This model is created for the so called variance stabilization property: asymptotically \hat{Y}_{ij} has variance only depending on number of ratings on edge $\{i, j\}$ or the weight ω_{ij} , but not on the true probability p_{ij} .

Different models will give different \hat{Y}_{ij} from the same observation $\hat{\pi}_{ij}$, followed by the same weighted least square problem (169) for the solution. Therefore, a deeper analysis of problem (169) will disclose more properties about the ranking problem.

HodgeRank on graph $G = (V, E)$ provides us such a tool, which characterizes the solution and residue of (169), adaptive to topological structures of G . The following theorem adapted from [JLYY11] describes a decomposition of \hat{Y} , which can be visualized as edge flows on graph G with direction $i \rightarrow j$ if $\hat{Y}_{ij} > 0$ and vice versa. Before the statement of the theorem, we first define the triangle set of G as all the 3-cliques in G .

$$(180) \quad T = \left\{ \{i, j, k\} \in \binom{V}{3} \mid \{i, j\}, \{j, k\}, \{k, i\} \in E \right\}.$$

Equipped with T , graph G becomes an abstract simplicial complex, the clique complex $\chi(G) = (V, E, T)$.

Theorem 1 [Hodge Decomposition of Paired Ranking] Let \hat{Y}_{ij} be a paired comparison flow on graph $G = (V, E)$, i.e., $\hat{Y}_{ij} = -\hat{Y}_{ji}$ for $\{i, j\} \in E$, and $\hat{Y}_{ij} = 0$ otherwise. There is a unique decomposition of \hat{Y} satisfying

$$(181) \quad \hat{Y} = \hat{Y}^g + \hat{Y}^h + \hat{Y}^c,$$

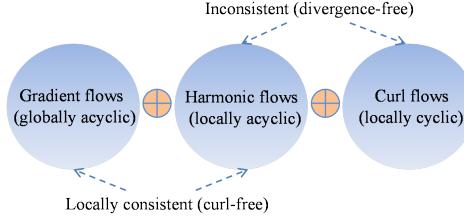


FIGURE 3. Hodge decomposition (three orthogonal components) of paired rankings [JLYY11].

where

$$(182) \quad \hat{Y}_{ij}^g = \hat{s}_i - \hat{s}_j, \text{ for some } \hat{s} \in \mathbb{R}^V,$$

$$(183) \quad \hat{Y}_{ij}^h + \hat{Y}_{jk}^h + \hat{Y}_{ki}^h = 0, \text{ for each } \{i, j, k\} \in T,$$

$$(184) \quad \sum_{j \sim i} \omega_{ij} \hat{Y}_{ij}^h = 0, \text{ for each } i \in V.$$

The decomposition above is *orthogonal* under the following inner product on $\mathbb{R}^{|E|}$, $\langle u, v \rangle_\omega = \sum_{\{i,j\} \in E} \omega_{ij} u_{ij} v_{ij}$.

The following provides some remarks on the decomposition.

1. When G is connected, \hat{Y}_{ij}^g is a rank two skew-symmetric matrix and gives a linear score function $\hat{s} \in \mathbb{R}^V$ up to translations. We thus call \hat{Y}^g a *gradient flow* since it is given by the difference (discrete gradient) of the score function \hat{s} on graph nodes,

$$(185) \quad \hat{Y}_{ij}^g = (\delta_0 \hat{s})(i, j) := \hat{s}_i - \hat{s}_j,$$

where $\delta_0 : \mathbb{R}^V \rightarrow \mathbb{R}^E$ is a finite difference operator (matrix) on G . \hat{s} can be chosen as any least square solution of (169), where we often choose the minimal norm solution,

$$(186) \quad \hat{s} = \Delta_0^\dagger \delta_0^* \hat{Y},$$

where $\delta_0^* = \delta_0^T W$ ($W = \text{diag}(\omega_{ij})$), $\Delta_0 = \delta_0^* \cdot \delta_0$ is the unnormalized graph Laplacian defined by $(\Delta_0)_{ii} = \sum_{j \sim i} \omega_{ij}$ and $(\Delta_0)_{ij} = -\omega_{ij}$, and $(\cdot)^\dagger$ is the Moore-Penrose (pseudo) inverse. On a complete and balanced graph, (186) is reduced to $\hat{s}_i = \frac{1}{n-1} \sum_{j \neq i} \hat{Y}_{ij}$, often called *Borda Count* as the earliest preference aggregation rule in social choice [JLYY11]. For expander graphs like regular graphs, graph Laplacian Δ_0 has small condition numbers and thus the global ranking is stable against noise on data.

2. \hat{Y}^h satisfies two conditions (183) and (184), which are called *curl-free* and *divergence-free* conditions respectively. The former requires the triangular trace of \hat{Y} to be zero, on every 3-clique in graph G ; while the later requires the total sum (inflow minus outflow) to be zero on each node of G . These two conditions characterize a linear subspace which is called *harmonic flows*.

3. The residue \hat{Y}^c actually satisfies (184) but not (183). In fact, it measures the amount of intrinsic (local) inconsistency in \hat{Y} characterized by the triangular

trace. We often call this component *curl flow*. In particular, the following relative curl,

$$(187) \quad \text{curl}_{ijk}^r = \frac{|\hat{Y}_{ij} + \hat{Y}_{jk} + \hat{Y}_{ki}|}{|\hat{Y}_{ij}| + |\hat{Y}_{jk}| + |\hat{Y}_{ki}|} = \frac{|\hat{Y}_{ij}^c + \hat{Y}_{jk}^c + \hat{Y}_{ki}^c|}{|\hat{Y}_{ij}| + |\hat{Y}_{jk}| + |\hat{Y}_{ki}|} \in [0, 1],$$

can be used to characterize triangular intransitivity; $\text{curl}_{ijk}^r = 1$ iff $\{i, j, k\}$ contains an intransitive triangle of \hat{Y} . Note that computing the percentage of $\text{curl}_{ijk}^r = 1$ is equivalent to calculating the Transitivity Satisfaction Rate (TSR) in complete graphs.

Figure 3 illustrates the Hodge decomposition for paired comparison flows and Algorithm 12 shows how to compute global ranking and other components. The readers may refer to [JLYY11] for the detail of theoretical development. Below we just make a few comments on the application of HodgeRank in our setting.

Algorithm 12: Procedure of Hodge decomposition in Matlab Pseudocodes

Input: A paired comparison hypergraph G provide by assessors.
Output: Global score \hat{s} , gradient flow \hat{Y}^g , curl flow \hat{Y}^c , and harmonic flow \hat{Y}^h .

- 1 **Initialization:**
- 2 \hat{Y} (a numEdge-vector consisting \hat{Y}_{ij} defined),
- 3 W (a numEdge-vector consisting ω_{ij}).
- 4 **Step 1:**
- 5 Compute δ_0, δ_1 ; // δ_0 = gradient, δ_1 = curl
- 6 $\delta_0^* = \delta_0^T * \text{diag}(W)$; // the conjugate of δ_0
- 7 $\Delta_0 = \delta_0^* * \delta_0$; // Unnormalized Graph Laplacian
- 8 $\text{div} = \delta_0^* * \hat{Y}$; // divergence operator
- 9 $\hat{s} = \text{lsqr}(\Delta_0, \text{div})$; // global score
- 10 **Step 2:**
- 11 Compute 1st projection on gradient flow: $\hat{Y}^g = \delta_0 * \hat{s}$;
- 12 **Step 3:**
- 13 $\delta_1^* = \delta_1^T * \text{diag}(1./W)$;
- 14 $\Delta_1 = \delta_1 * \delta_1^*$;
- 15 $\text{curl} = \delta_1 * \hat{Y}$;
- 16 $z = \text{lsqr}(\Delta_1, \text{curl})$;
- 17 Compute 3rd projection on curl flow: $\hat{Y}^c = \delta_1^* * z$;
- 18 **Step 4:**
- 19 Compute 2nd projection on harmonic flow: $\hat{Y}^h = \hat{Y} - \hat{Y}^g - \hat{Y}^c$.

1. To find a global ranking \hat{s} in (186), the recent developments of Spielman-Teng [ST04] and Koutis-Miller-Peng [KMP10] suggest fast (almost linear in $|E| \text{Poly}(\log |V|)$) algorithms for this purpose.

2. Inconsistency of \hat{Y} has two parts: global inconsistency measured by harmonic flow \hat{Y}^h and local inconsistency measured by curls in \hat{Y}^c . Due to the orthogonal decomposition, $\|\hat{Y}^h\|_\omega^2 / \|\hat{Y}\|_\omega^2$ and $\|\hat{Y}^c\|_\omega^2 / \|\hat{Y}\|_\omega^2$ provide percentages of global and local inconsistencies, respectively.

3. A nontrivial harmonic component $\hat{Y}^h \neq 0$ implies the fixed tournament issue, i.e., for any candidate $i \in V$, there is a paired comparison design by removing some of the edges in $G = (V, E)$ such that i is the overall winner.

4. One can control the harmonic component by controlling the topology of clique complex $\chi(G)$. In a loop-free clique complex $\chi(G)$ where $\beta_1 = 0$, harmonic component vanishes. In this case, there are no cycles which traverse all the nodes, *e.g.*, $1 \succ 2 \succ 3 \succ 4 \succ \dots \succ n \succ 1$. All the inconsistency will be summarized in those triangular cycles, *e.g.*, $i \succ j \succ k \succ i$.

Theorem 2. The linear space of harmonic flows has the dimension equal to β_1 , *i.e.*, the number of independent loops in clique complex $\chi(G)$, which is called the first order Betti number.

Fortunately, with the aid of some random sampling principles, it is not hard to obtain graphs whose β_1 are zero.

4.2. Random Graphs. In this section, we first describe two classical random models: Erdős-Rényi random graph and random regular graph; then we investigate the relation between them.

4.2.1. *Erdős-Rényi Random Graph.* Erdős-Rényi random graph $G(n, p)$ starts from n vertices and draws its edges independently according to a fixed probability p . Such random graph model is chosen to meet the scenario that in crowdsourcing ranking raters and videos come in an unspecified way. Among various models, Erdős-Rényi random graph is the simplest one equivalent to I.I.D. sampling. Therefore, such a model is to be systematically studied in the paper.

However, to exploit Erdős-Rényi random graph in crowdsourcing experimental designs, one has to meet some conditions depending on our purpose:

1. *The resultant graph should be connected, if we hope to derive global scores for all videos in comparison;*
2. *The resultant graph should be loop-free in its clique complex, if we hope to get rid of the global inconsistency in harmonic component.*

The two conditions can be easily satisfied for large Erdős-Rényi random graph.

Theorem 3. Let $G(n, p)$ be the set of Erdős-Rényi random graphs with n nodes and edge appearance probability p . Then the following holds as $n \rightarrow \infty$,

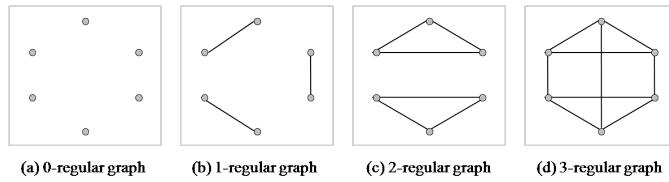
1. [Erdős-Rényi 1959] [ER59] if $p \succ \log n / n$, then $G(n, p)$ is almost always connected; and if $p \prec \log n / n$ then $G(n, p)$ is almost always disconnected;
2. [Kahle 2009] [Kah09, Kah13] if $p = O(n^\alpha)$, with $\alpha < -1$ or $\alpha > -1/2$, then the expected β_1 of the clique complex $\chi(G(n, p))$ is almost always equal to zero, *i.e.*, loop-free.

These theories imply that when p is large enough, Erdős-Rényi random graph will meet the two conditions above with high probability. In particular, almost linear $O(n \log n)$ edges suffice to derive a global ranking, and with $O(n^{3/2})$ edges harmonic-free condition is met.

Despite such an asymptotic theory for large random graphs, it remains a question how to ensure that a given graph instance satisfies the two conditions? Fortunately, the recent development in computational topology provides us such a tool, persistent homology, which will be illustrated in Section ??.

5. Euler-Calculus

to be finished...

FIGURE 4. Examples of k -regular graphs.

Methods	K-means	K-center	Average	Complete	Single
Complexity	NP	NP	\approx K-means	\approx K-center	Minimal Spanning Tree
Approximability	50-opt	2-opt. $O(kn)$?	$k < \alpha(k) < k^{\log 3}$?
Online	?	Cover-tree (8-opt)	?	?	Persistent Homology
Hierarchical	?	Cover-tree	Yes	Yes	Yes
Consistency	Pollard'81	No (metric-net)	?	?	Hartigen'81; Stuetzle'03

Bibliography

- [AC09] R DeVore A Cohen, W Dahmen, *Compressed sensing and best k -term approximation*, J. Amer. Math. Soc **22** (2009), no. 1, 211–231.
- [Ach03] Dimitris Achlioptas, *Database-friendly random projections: Johnson-lindenstrauss with binary coins*, Journal of Computer and System Sciences **66** (2003), 671687.
- [Ali95] F. Alizadeh, *Interior point methods in semidefinite programming with applications to combinatorial optimization*, SIAM J. Optim. **5** (1995), no. 1, 13–51.
- [Aro50] N. Aronszajn, *Theory of reproducing kernels*, Transactions of the American Mathematical Society **68** (1950), no. 3, 337–404.
- [Bav11] Francois Bavaud, *On the schoenberg transformations in data analysis: Theory and illustrations*, Journal of Classification **28** (2011), no. 3, 297–314.
- [BDDW08] Richard Baraniuk, Mark Davenport, Ronald DeVore, and Michael Wakin, *A simple proof of the restricted isometry property for random matrices*, Constructive Approximation **28** (2008), no. 3, 253–263.
- [BE92] Andreas Buja and Nermin Eyuboglu, *Remarks on parallel analysis*, Multivariate Behavioral Research **27** (1992), no. 4, 509–540.
- [BFOS07] M. Burger, K. Frick, S. Osher, and O. Scherzer, *Inverse total variation flow*, SIAM Multiscale Model. Simul. **6** (2007), no. 2, 366–395.
- [BGOX06] Martin Burger, Guy Gilboa, Stanley Osher, and Jinjun Xu, *Nonlinear inverse scale space methods*, Communications in Mathematical Sciences **4** (2006), no. 1, 179–212.
- [BLT⁺06] P. Biswas, T.-C. Liang, K.-C. Toh, T.-C. Wang, and Y. Ye, *Semidefinite programming approaches for sensor network localization with noisy distance measurements*, IEEE Transactions on Automation Science and Engineering **3** (2006), 360–371.
- [BN01] Mikhail Belkin and Partha Niyogi, *Laplacian eigenmaps and spectral techniques for embedding and clustering*, Advances in Neural Information Processing Systems (NIPS) 14, MIT Press, 2001, pp. 585–591.
- [BN03] Mikhail Belkin and Partha Niyogi, *Laplacian eigenmaps for dimensionality reduction and data representation*, Neural Computation **15** (2003), 1373–1396.
- [BN08] Mikhail Belkin and Partha Niyogi, *Convergence of laplacian eigenmaps*, Tech. report, 2008.
- [BP98] Sergey Brin and Larry Page, *The anatomy of a large-scale hypertextual web search engine*, Proceedings of the 7th international conference on World Wide Web (WWW) (Australia), 1998, pp. 107–117.
- [BS10] Zhidong Bai and Jack W. Silverstein, *Spectral analysis of large dimensional random matrices*, Springer, 2010.
- [BTA04] Alain Berlinet and Christine Thomas-Agnan, *Reproducing kernel hilbert spaces in probability and statistics*, Kluwer Academic Publishers, 2004.
- [Bur08] Martin Burger, *A note on sparse reconstruction methods*, Journal of Physics Conference Series **124** (2008), no. 1, 012002.
- [Can08] E. J. Candès, *The restricted isometry property and its implications for compressed sensing*, Comptes Rendus de l'Académie des Sciences, Paris, Série I **346** (2008), 589–592.
- [CDS98] Scott Shaobing Chen, David L. Donoho, and Michael A. Saunders, *Atomic decomposition by basis pursuit*, SIAM Journal on Scientific Computing **20** (1998), 33–61.
- [Chu05] Fan R. K. Chung, *Laplacians and the cheeger inequality for directed graphs*, Annals of Combinatorics **9** (2005), no. 1, 1–19.
- [CL06] Ronald R. Coifman and Stéphane. Lafon, *Diffusion maps*, Applied and Computational Harmonic Analysis **21** (2006), 5–30.

- [CLL⁺05] R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner, and S. W. Zucker, *Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps i*, Proceedings of the National Academy of Sciences of the United States of America **102** (2005), 7426–7431.
- [CLMW09] E. J. Candès, Xiaodong Li, Yi Ma, and John Wright, *Robust principal component analysis*, Journal of ACM **58** (2009), no. 1, 1–37.
- [CMOP11] Ozan Candogan, Ishai Menache, Asuman Ozdaglar, and Pablo A. Parrilo, *Flows and decompositions of games: Harmonic and potential games*, Mathematics of Operations Research **36** (2011), no. 3, 474–503.
- [Coo07] R. Dennis Cook, *Fisher lecture: Dimension reduction in regression*, Statistical Science **22** (2007), no. 1, 1–26.
- [CPW12] V. Chandrasekaran, P. A. Parrilo, and A. S. Willsky, *Latent variable graphical model selection via convex optimization (with discussion)*, Annals of Statistics (2012), to appear, <http://arxiv.org/abs/1008.1290>.
- [CR09] E. J. Candès and B. Recht, *Exact matrix completion via convex optimization*, Foundation of Computational Mathematics **9** (2009), no. 6, 717772.
- [CRPW12] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, *The convex geometry of linear inverse problems*, Foundation of Computational Mathematics (2012), to appear, <http://arxiv.org/abs/1012.0621>.
- [CRT06] Emmanuel J. Candès, Justin Romberg, and Terrence Tao, *Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information*, IEEE Trans. on Info. Theory **52** (2006), no. 2, 489–509.
- [CSPW11] V. Chandrasekaran, S. Sanghavi, P.A. Parrilo, and A. Willsky, *Rank-sparsity incoherence for matrix decomposition*, SIAM Journal on Optimization **21** (2011), no. 2, 572–596, <http://arxiv.org/abs/0906.2220>.
- [CST03] N. Cristianini and J. Shawe-Taylor, *An introduction to support vector machines and other kernel-based learning methods*, Cambridge University Press, 2003.
- [CT05] E. J. Candès and Terrence Tao, *Decoding by linear programming*, IEEE Trans. on Info. Theory **51** (2005), 4203–4215.
- [CT06] Emmanuel J. Candès and Terrence Tao, *Near optimal signal recovery from random projections: Universal encoding strategies*, IEEE Trans. on Info. Theory **52** (2006), no. 12, 5406–5425.
- [CT10] E. J. Candès and T. Tao, *The power of convex relaxation: Near-optimal matrix completion*, IEEE Transaction on Information Theory **56** (2010), no. 5, 2053–2080.
- [Dav88] H. David, *The methods of paired comparisons*, 2nd ed., Griffin's Statistical Monographs and Courses, 41, Oxford University Press, New York, NY, 1988.
- [DG03a] Sanjoy Dasgupta and Anupam Gupta, *An elementary proof of a theorem of johnson and lindenstrauss*, Random Structures and Algorithms **22** (2003), no. 1, 60–65.
- [DG03b] David L. Donoho and Carrie Grimes, *Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data*, Proceedings of the National Academy of Sciences of the United States of America **100** (2003), no. 10, 5591–5596.
- [dGJL07] Alexandre d'Aspremont, Laurent El Ghaoui, Michael I. Jordan, and Gert R. G. Lanckriet, *A direct formulation for sparse pca using semidefinite programming*, SIAM Review **49** (2007), no. 3, <http://arxiv.org/abs/cs/0406021>.
- [DH01] David L. Donoho and Xiaoming Huo, *Uncertainty principles and ideal atomic decomposition*, IEEE Transactions on Information Theory **47** (2001), no. 7, 2845–2862.
- [EB01] M. Elad and A.M. Bruckstein, *On sparse representations*, International Conference on Image Processing (ICIP) (Tsaloniky, Greece), November 2001.
- [Efr10] Bradley Efron, *Large-scale inference: Empirical bayes methods for estimation, testing, and prediction*, Cambridge University Press, 2010.
- [ELVE08] Weinan E, Tiejun Li, and Eric Vanden-Eijnden, *Optimal partition and effective dynamics of complex networks*, Proc. Nat. Acad. Sci. **105** (2008), 7907–7912.
- [ER59] P. Erdos and A. Renyi, *On random graphs i*, Publ. Math. Debrecen **6** (1959), 290–297.
- [EST09] Ioannis Z. Emiris, Frank J. Sottile, and Thorsten Theobald, *Nonlinear computational geometry*, Springer, New York, 2009.
- [EVE06] Weinan E and Eric Vanden-Eijnden, *Towards a theory of transition paths*, J. Stat. Phys. **123** (2006), 503–523.

- [EVE10] Weinan E and Eric Vanden-Eijnden, *Transition-path theory and path-finding algorithms for the study of rare events*, Annual Review of Physical Chemistry **61** (2010), 391–420.
- [FHX⁺16] Yanwei Fu, Timothy M. Hospedales, Tao Xiang, Jiechao Xiong, Shaogang Gong, Yizhou Wang, and Yuan Yao, *Robust subjective visual property prediction from crowdsourced pairwise labels*, IEEE Transactions on Pattern Analysis and Machine Intelligence **38** (2016), no. 3, 563–577.
- [FL01] Jianqing Fan and Runze Li, *Variable selection via nonconcave penalized likelihood and its oracle properties*, Journal of American Statistical Association (2001), 1348–1360.
- [Gro11] David Gross, *Recovering low-rank matrices from few coefficients in any basis*, IEEE Transaction on Information Theory **57** (2011), 1548, [arXiv:0910.1879](https://arxiv.org/abs/0910.1879).
- [HAvL05] M. Hein, J. Audibert, and U. von Luxburg, *From graphs to manifolds: weak and strong pointwise consistency of graph laplacians*, COLT, 2005.
- [Hor65] John L. Horn, *A rationale and test for the number of factors in factor analysis*, Psychometrika **30** (1965), no. 2, 179–185.
- [HS89] Trevor Hastie and Werner Stuetzle, *Principal curves*, Journal of the American Statistical Association **84** (1989), no. 406, 502–516.
- [HSXY16] Chendi Huang, Xinwei Sun, Jiechao Xiong, and Yuan Yao, *Split lbi: An iterative regularization path with structural sparsity*, Advances in Neural Information Processing Systems (NIPS) 29 (D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds.), 2016, pp. 3369–3377.
- [HSXY18] ———, *Boosting with structural sparsity: A differential inclusion approach*, Applied and Computational Harmonic Analysis (2018).
- [HTF01] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The elements of statistical learning*, Springer, 2001.
- [Hub81] P. J. Huber, *Robust statistics*, New York: Wiley, 1981.
- [HY18] Chendi Huang and Yuan Yao, *A unified dynamic approach to sparse model selection*, The 21st International Conference on Artificial Intelligence and Statistics (AISTATS) (Lanzarote, Spain), 2018.
- [JL84] W. B. Johnson and J. Lindenstrauss, *Extensions of lipschitz maps into a hilbert space*, Contemp Math **26** (1984), 189–206.
- [JLYY11] Xiaoye Jiang, Lek-Heng Lim, Yuan Yao, and Yinyu Ye, *Statistical ranking and combinatorial hodge theory*, Mathematical Programming **127** (2011), no. 1, 203–244, [arXiv:0811.1067](https://arxiv.org/abs/0811.1067) [stat.ML].
- [Joh06] I. Johnstone, *High dimensional statistical inference and random matrices*, Proc. International Congress of Mathematicians, 2006.
- [JYLG12] Xiaoye Jiang, Yuan Yao, Han Liu, and Leo Guibas, *Detecting network cliques with radon basis pursuit*, The Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS) (La Palma, Canary Islands), April 21-23 2012.
- [Kah09] Matthew Kahle, *Topology of random clique complexes*, Discrete Mathematics **309** (2009), 1658–1671.
- [Kah13] ———, *Sharp vanishing thresholds for cohomology of random flag complexes*, Annals of Mathematics (2013), arXiv:1207.0149.
- [Kle99] Jon Kleinberg, *Authoritative sources in a hyperlinked environment*, Journal of the ACM **46** (1999), no. 5, 604–632.
- [KMP10] Ioannis Koutis, G. Miller, and Richard Peng, *Approaching optimality for solving sdd systems*, FOCS ’10 51st Annual IEEE Symposium on Foundations of Computer Science, 2010, pp. 235–244.
- [KN08] S. Krutchman and B. Nadler, *Determining the number of components in a factor model from limited noisy data*, Chemometrics and Intelligent Laboratory Systems **94** (2008), 19–32.
- [Li91] Ker-Chau Li, *Sliced inverse regression for dimension reduction*, Journal of the American Statistical Association **86** (1991), no. 414, 316–327.
- [LK10] Dandan Li and Chung-Ping Kwong, *Understanding latent semantic indexing: A topological structure analysis using q-analysis*, J. Am. Soc. Inf. Sci. Technol. **61** (2010), no. 3, 592–608.

- [LL11] Jian Li and Tiejun Li, *Probabilistic framework for network partition*, Phys. A **390** (2011), 3579.
- [LLE09] Tiejun Li, Jian Liu, and Weinan E, *Probabilistic framework for network partition*, Phys. Rev. E **80** (2009), 026106.
- [LM06] Amy N. Langville and Carl D. Meyer, *Google's pagerank and beyond: The science of search engine rankings*, Princeton University Press, 2006.
- [LST13] Jason D Lee, Yuekai Sun, and Jonathan E Taylor, *On model selection consistency of penalized m -estimators: a geometric theory*, Advances in Neural Information Processing Systems (NIPS) 26, 2013, pp. 342–350.
- [LZ10] Yanhua Li and Zhili Zhang, *Random walks on digraphs, the generalized digraph laplacian, and the degree of asymmetry*, Algorithms and Models for the Web-Graph, Lecture Notes in Computer Science, vol. 6516, 2010, pp. 74–85.
- [LZ11] Gilad Lerman and Teng Zhang, *Robust recovery of multiple subspaces by geometric l_p minimization*, Annals of Statistics **39** (2011), no. 5, 2686–2715.
- [Mey00] Carl D. Meyer, *Matrix analysis and applied linear algebra*, SIAM, 2000.
- [MSVE09] Philipp Metzner, Christof Schütte, and Eric Vanden-Eijnden, *Transition path theory for markov jump processes*, Multiscale Model. Simul. **7** (2009), 1192.
- [MY09] Nicolai Meinshausen and Bin Yu, *Lasso-type recovery of sparse representations for high-dimensional data*, Annals of Statistics **37** (2009), no. 1, 246–270.
- [MZ93] S. G. Mallat and Z. Zhang, *Matching pursuits with time-frequency dictionaries*, IEEE Transactions on Signal Processing **41** (1993), no. 12, 3397–3415.
- [NBG10] R. R. Nadakuditi and F. Benaych-Georges, *The breakdown point of signal subspace estimation*, IEEE Sensor Array and Multichannel Signal Processing Workshop (2010), 177–180.
- [Noe60] G. Noether, *Remarks about a paired comparison model*, Psychometrika **25** (1960), 357–367.
- [NSVE⁺09] Frank Noè, Christof Schütte, Eric Vanden-Eijnden, Lothar Reich, and Thomas R. Weikl, *Constructing the equilibrium ensemble of folding pathways from short off-equilibrium simulations*, Proceedings of the National Academy of Sciences of the United States of America **106** (2009), no. 45, 19011–19016.
- [OBG⁺05] Stanley Osher, Martin Burger, Donald Goldfarb, Jinjun Xu, and Wotao Yin, *An iterative regularization method for total variation-based image restoration*, SIAM Journal on Multiscale Modeling and Simulation **4** (2005), no. 2, 460–489.
- [ORX⁺16] Stanley Osher, Feng Ruan, Jiechao Xiong, Yuan Yao, and Wotao Yin, *Sparse recovery via differential inclusions*, Applied and Computational Harmonic Analysis **41** (2016), no. 2, 436–469, [arXiv:1406.7728](https://arxiv.org/abs/1406.7728).
- [OW16] Art Owen and Jingshu Wang, *Bi-cross-validation for factor analysis*, Statist. Sci. **31** (2016), no. 1, 119–139.
- [RL00] Sam T. Roweis and Saul K. Lawrence, *Locally linear embedding*, Science **290** (2000), no. 5500, 2319–2323.
- [RXY18] Feng Ruan, Jiechao Xiong, and Yuan Yao, *Libra: Linearized bregman algorithms for generalized linear models*, 2018, R package version 1.6, <https://cran.r-project.org/web/packages/Libra>.
- [Sch37] I. J. Schoenberg, *On certain metric spaces arising from euclidean spaces by a change of metric and their imbedding in hilbert space*, The Annals of Mathematics **38** (1937), no. 4, 787–793.
- [Sch38a] ———, *Metric spaces and completely monotone functions*, The Annals of Mathematics **39** (1938), 811–841.
- [Sch38b] ———, *Metric spaces and positive definite functions*, Transactions of the American Mathematical Society **44** (1938), 522–536.
- [SHYW17] Xinwei Sun, Lingjing Hu, Yuan Yao, and Yizhou Wang, *Gsplit lbi: Taming the procedural bias in neuroimaging for disease prediction*, International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI), Springer, 2017, pp. 107–115.
- [Sin06] Amit Singer, *From graph to manifold laplacian: The convergence rate*, Applied and Computational Harmonic Analysis **21** (2006), 128–134.

- [ST04] D. Spielman and Shang-Hua Teng, *Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems*, STOC '04 Proceedings of the thirty-sixth annual ACM symposium on Theory of computing, 2004.
- [Ste56] Charles Stein, *Inadmissibility of the usual estimator for the mean of a multivariate distribution*, Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability **1** (1956), 197–206.
- [SW12] Amit Singer and Hau-Tieng Wu, *Vector diffusion maps and the connection laplacian*, Comm. Pure Appl. Math. **65** (2012), no. 8, 1067–1144.
- [SY07] Anthony Man-Cho So and Yinyu Ye, *Theory of semidefinite programming for sensor network localization*, Mathematical Programming, Series B **109** (2007), no. 2-3, 367–384.
- [SYZ08] Anthony Man-Cho So, Yinyu Ye, and Jiawei Zhang, *A unified theorem on sdp rank reduction*, Mathematics of Operations Research **33** (2008), no. 4, 910–920.
- [Tao11] Terrence Tao, *Topics in random matrix theory*, Lecture Notes in UCLA, 2011.
- [TdL00] J. B. Tenenbaum, Vin deSilva, and John C. Langford, *A global geometric framework for nonlinear dimensionality reduction*, Science **290** (2000), 2319–2323.
- [TdSL00] J. Tenenbaum, V. de Silva, and J. Langford, *A global geometric framework for nonlinear dimensionality reduction*, Science **290** (2000), no. 5500, 2323–2326.
- [Tib96] R. Tibshirani, *Regression shrinkage and selection via the lasso*, J. of the Royal Statistical Society, Series B **58** (1996), no. 1, 267–288.
- [Tro04] Joel A. Tropp, *Greed is good: Algorithmic results for sparse approximation*, IEEE Trans. Inform. Theory **50** (2004), no. 10, 2231–2242.
- [Tsy09] Alexandre Tsybakov, *Introduction to nonparametric estimation*, Springer, 2009.
- [Tyl87] D. E. Tyler, *A distribution-free m-estimator of multivariate scatter*, Annals of Statistics **15** (1987), no. 1, 234–251.
- [Vap98] V. Vapnik, *Statistical learning theory*, Wiley, New York, 1998.
- [Vem04] Santosh Vempala, *The random projection method*, Am. Math. Soc., Providence, 2004.
- [Wah90] Grace Wahba, *Spline models for observational data*, CBMS-NSF Regional Conference Series in Applied Mathematics 59, SIAM, 1990.
- [WLM09] Qiang Wu, Feng Liang, and Sayan Mukherjee, *Localized sliced inverse regression*, Annual Conference on Neural Information Processing Systems (NIPS) (2009).
- [WS06] Killian Q. Weinberger and Lawrence K. Saul, *Unsupervised learning of image manifolds by semidefinite programming*, International Journal of Computer Vision **70** (2006), no. 1, 77–90.
- [XRY18] Jiechao Xiong, Feng Ruan, and Yuan Yao, *A tutorial on libra: R package for the linearized bregman algorithms in high dimensional statistics*, Handbook of Big Data Analytics, Springer, 2018.
- [XXCY16a] Qianqian Xu, Jiechao Xiong, Xiaochun Cao, and Yuan Yao, *False discovery rate control and statistical quality assessment of annotators in crowdsourced ranking*, International Conference on Machine Learning (ICML), 2016, New York, June 19-24.
- [XXCY16b] ———, *Parsimonious mixed-effects HodgeRank for crowdsourced preference aggregation*, ACM Multimedia Conference, 2016.
- [YH41] G. Young and A. S. Householder, *A note on multidimensional psycho-physical analysis*, Psychometrika **6** (1941), 331–333.
- [YODG08] Wotao Yin, Stanley Osher, Jerome Darbon, and Donald Goldfarb, *Bregman iterative algorithms for compressed sensing and related problems*, SIAM Journal on Imaging Sciences **1** (2008), no. 1, 143–168.
- [ZCS14] Teng Zhang, Xiuyuan Cheng, and Amit Singer, *Marcenko-pastur law for tyler's m-estimator*.
- [Zha16] Teng Zhang, *Robust subspace recovery by tyler's m-estimator*, Information and Inference: A Journal of the IMA (2016), 1–23.
- [ZHT06] H. Zou, T. Hastie, and R. Tibshirani, *Sparse principal component analysis*, Journal of Computational and Graphical Statistics **15** (2006), no. 2, 262–286.
- [ZSF⁺18] Bo Zhao, Xinwei Sun, Yanwei Fu, Yuan Yao, and Yizhou Wang, *Msplit lbi: Realizing feature selection and dense estimation simultaneously in few-shot and zero-shot learning*, International Conference on Machine Learning (ICML), 2018.
- [ZW] Zhenyue Zhang and Jing Wang, *Mlle: Modified locally linear embedding using multiple weights*, <http://citeseex.ist.psu.edu/viewdoc/summary?doi=10.1.1.70.382>.

- [ZY06] Peng Zhao and Bin Yu, *On model selection consistency of lasso*, J. Machine Learning Research **7** (2006), 2541–2567.
- [ZZ02] Zhenyue Zhang and Hongyuan Zha, *Principal manifold and nonlinear dimension reduction via local tangent space alignment*, SIAM Journal of Scientific Computing **26** (2002), 313–338.
- [ZZ09] Hongyuan Zha and Zhenyue Zhang, *Spectral properties of the alignment matrices in manifold learning*, SIAM Review **51** (2009), no. 3, 545–566.