

國立臺灣大學理學院應用數學科學研究所

碩士論文

Institute of Applied Mathematical Sciences

College of Science

National Taiwan University

Master Thesis



整合多重隨機奇異值分解與理論分析

Theoretical and Performance Analysis for Integrated
Randomized Singular Value Decomposition

張大衛

Da-Wei Chang

指導教授：王偉仲博士

Advisor: Weichung Wang, Ph.D.

中華民國 106 年 7 月

July, 2017



國立臺灣大學碩士學位論文
口試委員會審定書



整合多重隨機奇異值分解與理論分析

Theoretical and Performance Analysis for Integrated
Randomized Singular Value Decomposition

本論文係張大衛君 (R04246002) 在國立臺灣大學應用數學科學研究所完成之碩士學位論文，於民國 106 年 7 月 26 日承下列考試委員審查通過及口試及格，特此證明

口試委員：

王信仲

陳定立

陳孝登

所長：





誌謝

在最先開始，我想感謝與這篇論文相關的所有人，沒有你們的協助，我無法將這篇論文完成。

感謝我的指導教授王偉仲老師，教導我該如何進行研究，定期開會討論進度，平常也總是開著辦公室的門，當遭遇問題，或是有新的進展時，都可以一起討論，並給予建設性的回應，除了學術討論以外，也對寫作與報告的模式進行指導，讓我對於論文的寫作有更深入的認識。感謝中研院統計所陳素雲老師與陳定立老師讓我參與整合隨機奇異值分解的這項研究，並且在理論方面提供指導與協助。感謝我的同學楊慕與林建耀，陪伴我一起討論，並且提出了許多新的點子與想法。感謝 R439 研究室的所有同學，營造出歡樂的研究氣氛，推動研究的推展，並且舒緩了研究不順利時的煩躁感。感謝我的家人在這段時間裡在生活與心理上的支持，使我可以無憂無慮，專心一致地完成這篇論文。

最後，再次感謝所有與這篇論文相關的所有人，從一開始接觸問題，到現在完成這份論文，因為有敬愛的師長、同學與家長的支持，才能讓這篇論文順利誕生。





Acknowledgements

First of all, I would like to thank my advisor Professor Weichung Wang of the Institute of Applied Mathematical Science at National Taiwan University. When I ran into a trouble of the research, Prof. Wang always opened his office door and discussed with me. He also spent a lot of time teaching me the skill for academic writing and presentation. Without his passionate guidance, this thesis would not have been successfully conducted.

I would also like to thank the collaborators of my advisor Dr. Su-Yun Huang and Dr. Ting-Li Chen of the Institute of Statistical Science at Academia Sinica. They allowed me to join this research and guided me in theoretical research. With their aid, the theory in this thesis could be vigorously constructed.

At the end, I would like to acknowledge my classmates, Mu Yang and Chienyao Lin. I am very grateful for their comments about this research. Some analysis and methods in this research are even inspired by their interesting ideas.





摘要

維度降低和特徵提取是大數據時代的重要技術，此二技術可以降低數據維數並降低進一步分析數據的計算成本。低秩奇異值分解 (low-rank SVD) 是這些技術的關鍵部分。為了更快地計算低秩奇異值分解，一些研究提出可以使用隨機抽取子空間的方法來獲得近似結果。在這項研究中，我們提出了一種新的概念，將隨機算法的結果進行整合以獲得更準確的近似值，稱為整合奇異值分解。我們通過理論和數值實驗來分析演算法的性質，以及不同的整合方法。整合方法的架構是有條件的優化問題，其具有唯一的局部極小值。整合子空間將透過線搜索、Kolmogorov-Nagumo 平均、和簡化類型的方法來進行計算，並針對這些方法的理論背景及計算複雜度進行分析，此外，整合奇異值分解與先前隨機奇異值分解的相似與相異處也會進行說明與分析。數值實驗結果顯示，在所提供的例子中，整合奇異值分解相對於同樣數量的隨機奇異值分解，使用線搜索方法時的疊代次數較少。另外，使用簡化類型的方法，來當作線搜索方法的初始值，可以減少收斂所需的疊代次數。

關鍵詞： 數值線性代數、奇異值分解、隨機演算法、數值優化、維度降低





Abstract

Dimension reduction and feature extraction are the important techniques in the big-data era to reduce the dimension of data and the computational cost for further data analysis. Low-rank singular value decomposition (low-rank SVD) is the key part of these techniques. In order to compute low-rank SVD faster, some researchers propose to use randomized subspace sketching algorithm to get an approximation result (rSVD). In this research, we propose an idea for integrating the results from randomized algorithm to get a more accurate approximation, which is called integrated singular value decomposition (iSVD). We analyze iSVD and the integration methods by theoretical analysis and numerical experiment. The integration scheme is a constraint optimization problem with unique local maximizer up to orthogonal transformation. Line search type method, Kolmogorov-Nagumo type average method and reduction type method are introduced and analyzed for their theoretical background and computational complexity. The similarity and difference between iSVD and rSVD with same sketching number are also explained and analyzed. The numerical experiment shows that the line search method in iSVD converges faster than the one in rSVD for our test examples. Also, using the integrated subspace from reduction as the initial value of line search method can reduce the iteration number to converge.

Keywords: Numerical Linear Algebra, Singular Value Decomposition, Randomized Algorithm, Numerical Optimization, Dimension Reduction





Contents

口試委員會審定書	iii
誌謝	v
Acknowledgements	vii
摘要	ix
Abstract	xi
1 Introduction	1
2 Overview of Integrated Singular Value Decomposition	5
3 Properties of Integrated Subspace	9
3.1 Solution of the Optimization Problem	10
3.2 Asymptotic Behavior of the Integrated Subspace	12
3.3 Uniqueness of Local Maximizer	13
4 Integration Method	21
4.1 Line Search Type Method	21
4.2 Kolmogorov-Nagumo-Type Average	29
4.3 Reduction-Type Average	33
5 Comparison of rSVD and iSVD	37

6 Numerical Experiment	41
6.1 Different Number of Sketched Subspaces	42
6.2 Comparison of KN and WY	44
6.3 Comparison of iSVD, rSVD and Reduction	47
7 Discussion and Conclusion	51
Bibliography	53





List of Figures

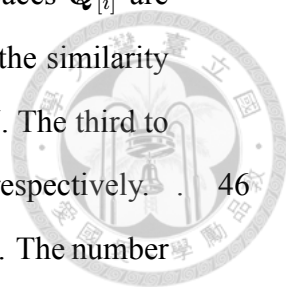
6.1 Similarity for different N . The size of test matrix is $m = 2^{19}$, $n = 2^{20}$. For each cases, we repeat 30 times iSVD with integration method WY and plot out the box plot of similarity. The box plot represent the maximum, Q3, median, Q1, minimum for each inner product among 30 times. 43

6.2 Average iteration number to converge for different N and different size of test matrix. The size of test matrix is $m = 2^d$, $n = 2^{d+1}$ for $d = 9, 11, 13, 15, 17, 19$. Each point shows the average iteration number among 30 tests. 43

6.3 Comparison of the approximate singular vectors by using WY and KN. The size of test matrix is $m = 2^{11}$, $n = 2^{12}$ and the sampling number $N = 32$. All of these test use the same 32 sketched subspaces $\mathbf{Q}_{[i]}$. The first line in the legend represents the similarity of $\overline{\mathbf{Q}} = \mathbf{Q}_{[1]}$. The second and third lines are the result from WY and KN respectively. The forth and fifth lines are from WY and KN respectively with fixed iteration number 15. 45

6.4 WY with different iteration numbers. The sketched subspaces $\mathbf{Q}_{[i]}$ are same in Figure 6.3. The first line in the legend represents the similarity for the case $\overline{\mathbf{Q}} = \mathbf{Q}_{[1]}$. The second line is the similarity from WY (with 61 iteration for $\mathbf{A}_H(10^{-1})$ and 84 iteration for $\mathbf{A}_H(10^{-3})$ to converge). The third to sixth lines are from WY with iteration number 5, 10, 15, 20 respectively. 46

6.5	KN with different iteration numbers. The sketched subspaces $Q_{[i]}$ are same in Figure 6.3. The first line in the legend represents the similarity for the case $Q_{[1]}$. The second line is the similarity from KN. The third to sixth lines are from KN with iteration number 5, 10, 15, 20 respectively.	46
6.6	Similarity for WY with iSVD and rSVD, reduction, and svds. The number of sketched subspaces in iSVD is $N = 32$. The number of sketching in rSVD is $32 * 22$, which is same as the total number of sketching in iSVD. The algorithm red.+WY uses the result of reduction as the initial value of WY.	48
6.7	Similarity for the methods in 6.6 with the fixed iteration number 10 for WY.	49





List of Tables

1.1	Notation in this thesis.	3
6.1	Abbreviation and detail information of the algorithm used in this section.	41





Chapter 1

Introduction

Dimension reduction and feature extraction are important issues in data analysis, especially for large scale data, to reduce the size or condense the information of analyzed data. The goal is to reduce the time for further analysis or find out the key features in the data. Singular value decomposition (SVD) is one of the technique to realize dimension reduction. An SVD of an $m \times n$ matrix \mathbf{A} takes the form $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where \mathbf{U} is an $m \times m$ orthogonal matrix, \mathbf{V} is an $n \times n$ orthogonal matrix, and $\mathbf{\Sigma}$ is an $m \times n$ diagonal matrix with decreasing diagonal entries. In this representation, \mathbf{U} , \mathbf{V} are the left and right singular vectors of \mathbf{A} respectively. The diagonal entries of $\mathbf{\Sigma}$ are the singular vector of \mathbf{A} . A rank- k approximation of \mathbf{A} via SVD is given as

$$\mathbf{A} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^\top$$

where \mathbf{U}_k , \mathbf{V}_k are leading k singular vectors and $\mathbf{\Sigma}_k$ contains leading k singular values. This low-rank approximation is called rank- k SVD of \mathbf{A} . Rank- k SVD is the best rank- k approximation of \mathbf{A} in the sense that it has the smallest 2-norm or Frobenius norm error. Therefore, it is a good choice for dimension reduction in many cases.

Many algorithms are aimed at computing the SVD or rank- k SVD. However, these algorithms usually take at least $O(m^2n + mnk)$ for computing rank- k SVD of a real $m \times n$ ($m \leq n$) matrix. Consequently, these algorithms cost lots of time to obtain the output result. Some research [7, 10] proposed a method to randomly sketch the matrix into a smaller

subspace, and calculate approximate rank- k SVD in that space. This method is called randomized singular value decomposition (rSVD) in this thesis. The accuracy and precision of rSVD depend on the quality of random projection used to generate the subspace. How to increase the quality of random sketching and the accuracy of the approximation is an important issue in rSVD. In [4], integrated singular value decomposition (iSVD) is proposed to enhance the projected subspace and obtain a better result. The main idea of iSVD is using integration method to condense the information from multiple randomly sketching and output a better-sketched subspace.

In this thesis, some properties of iSVD and integration method will be studied. The concept of iSVD is introduced in Chapter 2. The key step of iSVD is the integration, which is a constraint optimization problem. In the Chapter 3, we show the integrated subspace defined previously is the only local maximizer of the constraint optimization problem. The optimal solution and informal explanation of asymptotic behavior are introduced in the same chapter. These properties give the reason for using gradient type methods to solve this optimization problem. These methods are introduced and analyzed in Chapter 4. The similarity and difference between rSVD and iSVD with same sketching number are shown in the Chapter 5. The numerical experiment is shown in Chapter 6. Finally, the discussion and conclusion are given in Chapter 7.

The notations used in this thesis are as follows. The normal letters, such as a, α , denote the scalar. The bold lower case letters, such as $\mathbf{a}, \boldsymbol{\alpha}$, denote the vector. The bold upper case letters, such as $\mathbf{C}, \mathbf{\Gamma}$, denote the matrix. Table 1.1 shows some frequently appeared notations in this thesis.



\mathbf{A}	The matrix desired solving low-rank SVD.
m, n	The number of rows and columns of \mathbf{A} respectively. We assume $m \leq n$.
k	The given desired rank for low-rank approximation.
p	The number of oversampling in rSVD and iSVD.
ℓ	The total number of sampling in rSVD and iSVD for a single sketched subspaces. $\ell = k + p$. $\ell \ll m$.
N	The total number of sketched subspaces in iSVD.
\otimes	The Kronecker product.
$\mathbf{H}_{a,b}$	The $a \times b$ commutation matrix.[9] For any $a \times b$ matrix \mathbf{M} , $\text{vec}(\mathbf{M}^\top) = \mathbf{H}_{a,b} \text{vec}(\mathbf{M})$.
\mathbf{Q}_c	The current iterator in the iterative method.
\mathbf{Q}_+	The iterator of next step in the iterative method.

Table 1.1: Notation in this thesis.





Chapter 2

Overview of Integrated Singular Value Decomposition

Before introducing iSVD, we shall take a quick view of rSVD. The basic algorithm of rSVD from [7, 10] is stated as Algorithm 1. In this algorithm, SVD is only applied to an $m \times \ell$ matrix and an $\ell \times n$ matrix, which is much cheaper than applied on an $m \times n$ matrix if ℓ is small. The two main phases of rSVD are described in next two paragraphs.

Algorithm 1 Randomized SVD (rSVD)

Require: \mathbf{A} (real $m \times n$ matrix), k (number of desired rank for low-rank approximation), p (number of oversampling), $\ell = k + p$ (number of the sketched column),

Ensure: Approximate rank- k SVD of $\mathbf{A} \approx \hat{\mathbf{U}}_k \hat{\Sigma}_k \hat{\mathbf{V}}_k^\top$

- 1: Generate a random matrix $\mathbf{\Omega}$
 - 2: Assign $\mathbf{Y} \leftarrow \mathbf{A}\mathbf{\Omega}$
 - 3: Compute \mathbf{Q} whose columns are an orthonormal basis of \mathbf{Y}
 - 4: Compute the SVD of $\mathbf{Q}^\top \mathbf{A} = \widehat{\mathbf{W}}_\ell \widehat{\Sigma}_\ell \widehat{\mathbf{V}}_\ell^\top$
 - 5: Assign $\widehat{\mathbf{U}}_\ell \leftarrow \mathbf{Q}\widehat{\mathbf{W}}_\ell$
 - 6: Extract the largest k singular pairs from $\widehat{\mathbf{U}}_\ell, \widehat{\Sigma}_\ell, \widehat{\mathbf{V}}_\ell$ to obtain $\widehat{\mathbf{U}}_k, \widehat{\Sigma}_k, \widehat{\mathbf{V}}_k$
-

The first phase is randomly sketching a subspace of \mathbf{A} . More precisely, compute the matrix $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ and find out the orthogonal basis of the range of \mathbf{Y} as the approximate subspace \mathbf{Q} . In this case, $\mathbf{\Omega}$ is a randomly generated matrix. Both [7, 10] propose a commonly used $\mathbf{\Omega}$ as Gaussian projection, which means the entries of $\mathbf{\Omega}$ are independent identical standard normal distribution. This choice gives \mathbf{Y} some structure of the column space from \mathbf{A} and controls the error $\|\mathbf{Q}\mathbf{Q}^\top \mathbf{A} - \mathbf{A}\|_2 < \epsilon$ in high probability. Note that

$QQ^T A$ is the approximate rank- ℓ SVD of A with the column space spanned by Y .

The second phase is constructing an approximate rank- ℓ SVD $QQ^T A$. Only the SVD of $Q^T A$ is needed to construct the SVD of the rank- ℓ matrix $QQ^T A$. Once the SVD of $Q^T A = \widehat{W}_\ell \widehat{\Sigma}_\ell \widehat{V}_\ell^T$ is obtained, the SVD of $QQ^T A$ can be computed as $QQ^T A = (Q\widehat{W}_\ell) \widehat{\Sigma}_\ell \widehat{V}_\ell^T$. Note that the column space of this approximation is already determined in the first phase. The purpose of this phase is revealing the singular values and singular vectors in correct arrangement, and then the leading k singular values and singular vectors can be extracted.

The technique of oversampling is proposed here. Suppose the desired rank is k , the number of sketches in the first phase can be chosen as $\ell = k + p$ for some positive integer p . The approximate rank- k SVD is obtained from rank- ℓ SVD by extracting the first k singular vectors and singular values.

Based on rSVD, iSVD uses multiple random sketched subspaces in the first phase of rSVD to gain more accurate result. The algorithm is stated as Algorithm 2. Three phases are included in this algorithm. The first phase is similar to the first phase of SVD. Instead of choosing only one random sketch in rSVD, iSVD choose multiple random sketches. The second phase is integrating the subspaces obtained in the first phase and get an integrated subspace \overline{Q} . The third phase is same as the second phase of rSVD. They both construct the approximate rank- ℓ SVD.

Algorithm 2 Integrated SVD with multiple sketches (iSVD).

Require: A (real $m \times n$ matrix), k (desired rank of approximate SVD), p (oversampling parameter), $\ell = k + p$ (dimension of the sketched column space), q (exponent of the power method), N (number of random sketches)

Ensure: Approximate rank- k SVD of $A \approx \widehat{U}_k \widehat{\Sigma}_k \widehat{V}_k^T$

- 1: Generate $n \times \ell$ random matrices $\Omega_{[i]}$ for $i = 1, \dots, N$
 - 2: Assign $Y_{[i]} \leftarrow A\Omega_{[i]}$ for $i = 1, \dots, N$
 - 3: Compute $Q_{[i]}$ whose columns are an orthonormal basis of $Y_{[i]}$
 - 4: Integrate $\overline{Q} \leftarrow \{Q_{[i]}\}_{i=1}^N$
 - 5: Compute the SVD of $\overline{Q}^T A = \widehat{W}_\ell \widehat{\Sigma}_\ell \widehat{V}_\ell^T$
 - 6: Assign $\widehat{U}_\ell \leftarrow \overline{Q}\widehat{W}_\ell$
 - 7: Extract the largest k singular pairs from $\widehat{U}_\ell, \widehat{\Sigma}_\ell, \widehat{V}_\ell$ to obtain $\widehat{U}_k, \widehat{\Sigma}_k, \widehat{V}_k$
-

The key part of iSVD is how to define the integrated subspace in the second phase. The most intuitive idea is taking the arithmetic average of $Q_{[i]}$ as \overline{Q} , but the following state-

ments show this is not a reasonable definition of the integrated subspace. The integrated subspace \bar{Q} should represent the subspace integrated by $Q_{[i]}$. Each $Q_{[i]}$ is an orthogonal matrix. Hence \bar{Q} should also be an orthogonal matrix. However, the arithmetic average of $Q_{[i]}$ is not an orthogonal matrix. Therefore, the integrated subspaces should be defined by another form instead of the arithmetic average.

The integrated subspace in the second phase is defined by the following optimization problem.

$$\begin{cases} \bar{Q} := \operatorname{argmin}_Q \frac{1}{N} \sum_{i=1}^N \|Q_{[i]}Q_{[i]}^\top - QQ^\top\|_F^2 \\ \text{subject to } Q^\top Q = I_\ell \end{cases} \quad (2.1)$$

The main idea of this definition is to define the integrated subspace that has the minimum summation of distances between each $Q_{[i]}$ and \bar{Q} , which is similar to the property of the arithmetic average in Euclidean space. Instead of the Euclidean distance between Q and each $Q_{[i]}$, we use the distance between QQ^\top and $Q_{[i]}Q_{[i]}^\top$ to preserve the invariant of right orthogonal transformation. Suppose $Q_{[1]} = Q_{[2]}R$ for some orthogonal matrix R . $Q_{[1]}$ and $Q_{[2]}$ represent the same subspaces, so the measurement of error between $Q_{[1]}$, $Q_{[2]}$, and $Q_{[i]}$, $Q_{[i]}$ should be same. By using the relation

$$Q_{[1]}Q_{[1]}^\top - Q_{[i]}Q_{[i]}^\top = Q_{[1]}RR^\top Q_{[1]}^\top - Q_{[i]}Q_{[i]}^\top = Q_{[2]}Q_{[2]}^\top - Q_{[i]}Q_{[i]}^\top$$

the error measured in (2.1) gives the same value, which means it preserved the invariant of right orthogonal transformation. However, it still needs some theoretical analysis to make sure whether this objective function is suitable for the integration step. Chapter 3 provides an informal explanation that the integrated subspace defined by 2.1 can capture the true singular vectors of the desired matrix A when N tends to be large.





Chapter 3

Properties of Integrated Subspace

In this chapter, some properties of the integrated subspace defined in (2.1) are analyzed. First, an equivalent maximization problem (3.1) is introduced. Next, we show the solution of the optimization problem in (3.1) is the leading singular vectors of the arithmetic average of $\mathbf{Q}_{[i]}\mathbf{Q}_{[i]}^\top$. Then we provide an informal explanation for the asymptotic behavior of integrated subspace when N is large. Finally, we prove that the local maximizer of the optimization problem in (3.1) is unique up to a right orthogonal transform.

At the beginning, some equivalent problems should be introduced for the simplicity of analysis. The following theorem demonstrates an equivalent optimization problems of (2.1).

Theorem 3.0.1. *The constraint optimization problem (2.1) is equivalent to the problem*

$$\begin{cases} \operatorname{argmax}_{\mathbf{Q}} \frac{1}{2} \operatorname{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_\ell \end{cases} \quad (3.1)$$

where $\bar{\mathbf{P}} = \frac{1}{N} \sum_{i=1}^N \mathbf{Q}_{[i]}\mathbf{Q}_{[i]}^\top$ is the arithmetic average of $\mathbf{Q}_{[i]}\mathbf{Q}_{[i]}^\top$ over all i .

Proof. From the relation between trace and Frobenius norm $\|\mathbf{M}\|_F^2 = \operatorname{tr}(\mathbf{M}^\top \mathbf{M})$, the properties of trace $\operatorname{tr}(\mathbf{M}_1 \mathbf{M}_2) = \operatorname{tr}(\mathbf{M}_2 \mathbf{M}_1)$, and the orthogonality of \mathbf{Q} and $\mathbf{Q}_{[i]}$, the

Frobenius norm in the objective function of (2.1) can be rewritten as

$$\begin{aligned}
\|Q_{[i]}Q_{[i]}^\top - QQ^\top\|_F^2 &= \text{tr}((Q_{[i]}Q_{[i]}^\top - QQ^\top)^\top(Q_{[i]}Q_{[i]}^\top - QQ^\top)) \\
&= \text{tr}(QQ^\top) - 2\text{tr}(QQ^\top Q_{[i]}Q_{[i]}^\top) + \text{tr}(Q_{[i]}Q_{[i]}^\top) \\
&= \text{tr}(Q^\top Q) - 2\text{tr}(Q^\top Q_{[i]}Q_{[i]}^\top Q) + \text{tr}(Q_{[i]}^\top Q_{[i]}) \\
&= 2\ell - 2\text{tr}(Q^\top Q_{[i]}Q_{[i]}^\top Q).
\end{aligned}$$

This equation leads to another representation of the objective function

$$\frac{1}{N} \sum_{i=1}^N \|Q_{[i]}Q_{[i]}^\top - QQ^\top\|_F^2 = 2\frac{\ell}{N} - 2\frac{1}{N} \sum_{i=1}^N \text{tr}(Q^\top Q_{[i]}Q_{[i]}^\top Q) = 2\frac{\ell}{N} - 2\text{tr}(Q^\top \bar{P}Q).$$

Since the constant does not affect the problem and the negative scalar changes minimize problem to maximize problem, (2.1) is equivalent to the problem

$$\begin{cases} \operatorname{argmax}_Q \frac{1}{2} \text{tr}(Q^\top \bar{P}Q) \\ Q^\top Q = I_\ell \end{cases}$$

which is the problem (3.1). □

The optimization problem in the form (3.1) is more simple than the original form (2.1) for computing the derivative and further analyzing. So the simplified form (3.1) will be used in the following analysis.

3.1 Solution of the Optimization Problem

The solution of the maximize problem (3.1) is the integrated subspace we defined in (2.1). Theorem 3.1.1 shows that the optimizer of this problem is consisted by the leading ℓ eigenvectors of \bar{P} .

Theorem 3.1.1 (Maximal Value of Objective Function). *Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_\ell > \lambda_{\ell+1} \geq \dots \geq \lambda_m$ be the eigenvalues of \bar{P} and let Q_* be the corresponding leading ℓ*

eigenvectors. Then the objective function in (3.1) has the upper bound

$$\frac{1}{2} \text{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \leq \frac{1}{2} \sum_{i=1}^{\ell} \lambda_i$$

and the equality holds when $\mathbf{Q} = \mathbf{Q}_* \mathbf{R}_*$ for some $\ell \times \ell$ orthogonal matrix \mathbf{R}_* .



Proof. The eigenvalue decomposition of $\bar{\mathbf{P}}$ can be written as

$$\bar{\mathbf{P}} = \mathbf{Q}_* \mathbf{S}_* \mathbf{Q}_*^\top + \mathbf{Q}_\perp \mathbf{S}_\perp \mathbf{Q}_\perp^\top$$

where \mathbf{Q}_\perp denote orthonormal basis of the space perpendicular to \mathbf{Q} , $\mathbf{S}_* = \text{diag}(\lambda_1, \dots, \lambda_\ell)$ and $\mathbf{S}_\perp = \text{diag}(\lambda_{\ell+1}, \dots, \lambda_m)$. Since \mathbf{Q}_* and \mathbf{Q}_\perp spans the whole \mathbf{R}^m , any orthogonal matrix \mathbf{Q} can be represented as $\mathbf{Q} = \mathbf{Q}_* \mathbf{B} + \mathbf{Q}_\perp \mathbf{C}$ with $\mathbf{B}^\top \mathbf{B} + \mathbf{C}^\top \mathbf{C} = \mathbf{I}_\ell$. The objective function becomes

$$\begin{aligned} & \frac{1}{2} \text{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \\ &= \frac{1}{2} \text{tr}((\mathbf{Q}_* \mathbf{B} + \mathbf{Q}_\perp \mathbf{C})^\top (\mathbf{Q}_* \mathbf{S}_* \mathbf{Q}_*^\top + \mathbf{Q}_\perp \mathbf{S}_\perp \mathbf{Q}_\perp^\top) (\mathbf{Q}_* \mathbf{B} + \mathbf{Q}_\perp \mathbf{C})) \\ &= \frac{1}{2} \text{tr}(\mathbf{B}^\top \mathbf{S}_* \mathbf{B}) + \text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) \end{aligned}$$

by the relation $\mathbf{Q}^\top \mathbf{Q}_\perp = 0$

Let $\mathbf{B} = \mathbf{R} \mathbf{S} \mathbf{T}^\top$ be the SVD of \mathbf{B} . By multiplying $\mathbf{R} \mathbf{T}^\top$ from left and $\mathbf{T} \mathbf{R}^\top$ from right to the both sides of the condition $\mathbf{B}^\top \mathbf{B} + \mathbf{C}^\top \mathbf{C} = \mathbf{I}_\ell$, the new condition is given as $\mathbf{R} \mathbf{S}^2 \mathbf{R}^\top + \mathbf{R} \mathbf{T}^\top \mathbf{C}^\top \mathbf{C} \mathbf{T} \mathbf{R}^\top = \mathbf{R} \mathbf{T}^\top \mathbf{T} \mathbf{R}^\top = \mathbf{I}_\ell$. By using this new condition and the inequality

$$\text{tr}((\mathbf{C} \mathbf{T} \mathbf{R}^\top)^\top \mathbf{S}_\perp \mathbf{C} \mathbf{T} \mathbf{R}^\top) = \sum_{i=1}^{m-\ell} \sum_{j=1}^{\ell} \lambda_{\ell+i} c_{ij}^2 \leq \sum_{i=1}^{m-\ell} \sum_{j=1}^{\ell} \lambda_j c_{ij}^2 = \text{tr}(\mathbf{S}_* \mathbf{R} \mathbf{T}^\top \mathbf{C}^\top \mathbf{C} \mathbf{T} \mathbf{R}^\top)$$

where $CTR^\top = [c_{ij}]$, the upper bound of the objective function can be given as

$$\begin{aligned}
\frac{1}{2} \text{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) &= \frac{1}{2} \text{tr}(\mathbf{B}^\top \mathbf{S}_* \mathbf{B}) + \text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) \\
&= \frac{1}{2} (\text{tr}(\mathbf{T} \mathbf{S} \mathbf{R}^\top \mathbf{S}_* \mathbf{R} \mathbf{S} \mathbf{T}^\top) + \text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C} \mathbf{T} \mathbf{R}^\top \mathbf{R} \mathbf{T}^\top)) \\
&= \frac{1}{2} (\text{tr}(\mathbf{S}_* \mathbf{R} \mathbf{S}^2 \mathbf{R}^\top) + \text{tr}(\mathbf{R} \mathbf{T}^\top \mathbf{C}^\top \mathbf{S}_\perp \mathbf{C} \mathbf{T} \mathbf{R}^\top)) \\
&\leq \frac{1}{2} \text{tr}(\mathbf{S}_* (\mathbf{R} \mathbf{S}^2 \mathbf{R}^\top + \mathbf{R} \mathbf{T}^\top \mathbf{C}^\top \mathbf{C} \mathbf{T} \mathbf{R}^\top)) = \frac{1}{2} \text{tr}(\mathbf{S}_*) = \frac{1}{2} \sum_{i=1}^{\ell} \lambda_i.
\end{aligned}$$



The equality holds when $\sum_{i=1}^{m-\ell} \sum_{j=1}^{\ell} \lambda_{\ell+i} c_{ij}^2 = \sum_{i=1}^{m-\ell} \sum_{j=1}^{\ell} \lambda_j c_{ij}^2$. This equation means $CTR^\top = [c_{ij}] = 0$, and hence $\mathbf{C} = 0$. Therefore, $\mathbf{Q} = \mathbf{Q}_* \mathbf{R}_*$ with $\mathbf{B}^\top \mathbf{B} = \mathbf{I}_\ell$ and $\mathbf{R}_* = \mathbf{B}$. \square

Theorem 3.1.1 shows that the integrated subspace defined in (2.1) or (3.1) is actually formed by the leading eigenvectors of the arithmetic average of $\mathbf{Q}_{[i]} \mathbf{Q}_{[i]}^\top$. The leading eigenvectors of $\bar{\mathbf{P}}$ is same as the leading ℓ singular vectors of the matrix $[\mathbf{Q}_{[1]} | \mathbf{Q}_{[2]} | \cdots | \mathbf{Q}_{[N]}]$. This form is used for explaining the similarity of rSVD and iSVD with same sketching number in Chapter 5.

3.2 Asymptotic Behavior of the Integrated Subspace

Now we give an informal explanation about the reason why the the integrated subspace defined in (2.1) or (3.1) can work in the iSVD algorithm. Please refer to [4] for more detailed statistical analysis.

To explain the reason, the following theorem in [4] should be introduced first.

Theorem 3.2.1. *Let the SVD of \mathbf{A} be $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ with distinct decreasing singular values. Let \mathbf{Q} denote an orthogonal subspace spanned by $\mathbf{Y} = \mathbf{A} \mathbf{\Omega}$, where $\mathbf{\Omega}$ is randomly generated by i.i.d. standard normal entries. Then*

$$E [\mathbf{Q}_{[i]} \mathbf{Q}_{[i]}^\top] = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$$

where $\mathbf{\Lambda}$ satisfies

1. Λ is diagonal matrix

2. all diagonal entries are in the interval $(0, 1)$

3. the diagonal entries are decreasing if the singular values of \mathbf{A} are distinct.

By the law of large number and Theorem 3.2.1, as the sample size N goes large, the arithmetic average $\bar{\mathbf{P}}$ tends to a matrix with the same singular vectors of \mathbf{A} in the correct arrangement. Also, Theorem 3.1.1 shows that the optimizer defined in (3.1) is the leading singular vectors of $\bar{\mathbf{P}}$. Hence in the ideal case (N tends to infinity), iSVD can capture the leading singular vectors of \mathbf{A} as $\bar{\mathbf{Q}}$ and leads to an ideal result for the low-rank approximation.

3.3 Uniqueness of Local Maximizer

Recall that we focus on the problem in the form (3.1)

$$(+)\begin{cases} \max_{\mathbf{Q}} \frac{1}{2} \text{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_\ell \end{cases}$$

Theorem 3.1.1 shows the optimal solution of the problem (+) is formed by the leading ℓ eigenvectors of the matrix $\bar{\mathbf{P}}$. However, this result does not provide whether there exists another local maximum of this problem. The goal of the following derivation is Theorem 3.3.5, which shows that the only local maximum of the problem (+) is the orthogonal matrices that formed by the leading ℓ eigenvectors of $\bar{\mathbf{P}}$ (up to right orthogonal transformation). The main idea of the following proof is checking the first and second order necessary condition for the nonlinear equality constraint optimization problem. More information for the first and second order condition of optimization with equality constraint can be obtained in the book for introducing optimization problem. (For example, [6].)

We begin with checking the first order necessary condition of the problem (+).

Lemma 3.3.1 (First Order Necessary Condition). *Suppose \mathbf{Q} is a local maximizer of the problem (+) in the feasible set, i.e., the set collects all the $m \times \ell$ orthogonal matrix. Then*

Q satisfies the equation

$$(I - QQ^\top)\bar{P}Q = 0. \quad (3.2)$$

and the corresponding Lagrange multiplier $\Lambda \in \mathbf{R}^{\ell \times \ell}$ satisfies

$$\Lambda + \Lambda^\top = Q^\top \bar{P}Q. \quad (3.3)$$



Proof. This equation is obtained from the first order condition for optimal solution of an equality constraint, which states that if x^* is a local maximizer of an equality constraint optimization problem, then the following equations hold

$$\begin{cases} \nabla_x \mathcal{L}(x^*, \lambda^*) = 0 \\ \nabla_\lambda \mathcal{L}(x^*, \lambda^*) = 0 \end{cases}$$

where \mathcal{L} is the Lagrangian of the equality constraint optimization problem and λ^* is the corresponding Lagrange multiplier.

Notice that any $m \times n$ matrix M can be seen as an mn vector by vectorizing M as $\text{vec}(M)$. Also, one of the relation between $\text{vec}(\bullet)$ and Kronecker product \otimes is $\text{vec}(AXB) = (B^\top \otimes A) \text{vec}(X)$. By using these technique, the Lagrangian of (+) is given as

$$\begin{aligned} \mathcal{L}(Q, \Lambda) &= \frac{1}{2} \text{tr}(Q^\top \bar{P}Q) - \text{tr}(\Lambda^\top (Q^\top Q - I_\ell)) \\ &= \frac{1}{2} \text{vec}(Q)^\top \text{vec}(\bar{P}Q) - \text{vec}(\Lambda)^\top \text{vec}(Q^\top Q - I_\ell) \\ &= \frac{1}{2} \text{vec}(Q)^\top (I_\ell \otimes \bar{P}) \text{vec}(Q) - \text{vec}(\Lambda)^\top \text{vec}(Q^\top Q - I_\ell) \end{aligned}$$

and the first order conditions can be represented by the derivative of Lagrangian as

$$\begin{cases} \nabla_{\text{vec}(Q)} \mathcal{L}(Q, \Lambda) = (I \otimes \bar{P}) \text{vec}(Q) - [K_{\ell, m}(Q \otimes I_\ell) + (I_\ell \otimes Q)] \text{vec}(\Lambda) = 0 \\ \nabla_{\text{vec}(\Lambda)} \mathcal{L}(Q, \Lambda) = \text{vec}(Q^\top Q - I_\ell) = 0 \end{cases}$$

where $K_{\ell, m}$ denotes the commutation matrix corresponding to $\ell \times m$ matrix. By folding

the vectorized matrix back, the above equations leads to

$$\begin{cases} \overline{P}Q - Q(\Lambda + \Lambda^\top) = 0 \\ Q^\top Q - I_\ell = 0. \end{cases}$$



Applying Q^\top to the left of each side in the first equation and Using the second equation derives the relation $\Lambda + \Lambda^\top = Q^\top \overline{P}Q$. By substituting $(\Lambda + \Lambda^\top)$ as $Q^\top \overline{P}Q$ in the first equation, the first order conditions can be rewritten as

$$\begin{cases} (I - QQ^\top)\overline{P}Q = 0 \\ Q^\top Q = I_\ell \end{cases}$$

□

Lemma 3.3 gives the necessary conditions for feasible points and the corresponding Lagrange multiplier. The following lemma improves the result. It provides the explicit solutions that satisfy these conditions.

Lemma 3.3.2. *Let Q be a feasible points of the problem (+). Then Q satisfies the first order condition (3.2) if and only if each column of QR is an eigenvector of \overline{P} for some orthogonal matrix $R \in \mathbf{R}^{\ell \times \ell}$.*

Proof. For the ‘if’ part, since the columns of QR are eigenvectors of \overline{P} , the connection between \overline{P} and QR is given as $\overline{P}QR = QRS$, where S is a diagonal matrix with the corresponding eigenvalues on its diagonal entries. By directly calculation,

$$(I - QQ^\top)\overline{P}QR = (I - QQ^\top)QRS = 0$$

and hence

$$(I - QQ^\top)\overline{P}Q = 0$$

by multiplying R^\top to the right of both sides. This equation shows that Q satisfies the first order condition (3.2).

For the ‘only if’ part, suppose Q satisfies the first order condition $(I - QQ^\top)\bar{P}Q = 0$. Since $(I - QQ^\top)$ is the projection matrix onto the orthogonal complement of the range of Q , the first order condition means that the component of each columns of $\bar{P}Q$ in the orthogonal complement is 0 , i.e., each column of $\bar{P}Q$ is a vector in the range of Q . Therefore, each column of $\bar{P}Q$ can be written as the linear combination of the columns of Q and hence

$$\bar{P}Q = QB$$

for some $B \in \mathbf{R}^{\ell \times \ell}$. By applying Q^\top to the left of both sides, the equation becomes $B = Q^\top \bar{P}Q$ which is a symmetric matrix. Therefore, B has eigenvalue decomposition $B = RSR^\top$ for some orthogonal matrix R and diagonal matrix S . This decomposition gives the relation

$$\bar{P}Q = QRSR^\top$$

and hence

$$\bar{P}(QR) = (QR)S$$

by multiplying R to the right of both sides. This equation shows each column of (QR) is an eigenvector of \bar{P} since S is diagonal. \square

Now we check the second order condition. The preparation for computing second order condition is finding the null space spanned by the gradient of constrains. Let $g(Q)$ denote the constraint function $g(Q) = \text{vec}(Q^\top Q - I_\ell) = 0$, and \mathcal{N}_Q denote the null space of the linear space spanned by the gradients of constrains $\nabla_{\text{vec}(Q)}g(Q)$ at Q . The following lemma gives some explicit representations of the null space \mathcal{N}_Q .

Lemma 3.3.3. *The null space \mathcal{N}_Q is given as*

$$\mathcal{N}_Q = \{\text{vec}(Z) : Z \in \mathbf{R}^{m \times \ell}, Q^\top Z + Z^\top Q = 0\}$$

or equivalently,

$$\mathcal{N}_Q = \{\text{vec}(QB + Q_\perp C) : B \in \mathbf{R}^{\ell \times \ell}, C \in \mathbf{R}^{(m-\ell) \times \ell}, B = -B^\top\}$$

where \mathbf{Q}_\perp is a $m \times (m - \ell)$ orthogonal matrix satisfies $\mathbf{Q}_\perp^\top \mathbf{Q} = \mathbf{0}$, i.e., \mathbf{Q}_\perp contains the basis of orthogonal complements of the range of \mathbf{Q} .



Proof. By direct calculation, the gradient matrix of constraints at \mathbf{Q} is given as

$$\nabla_{\text{vec}(\mathbf{Q})} g(\mathbf{Q}) = \mathbf{K}_{\ell, m}(\mathbf{Q} \otimes \mathbf{I}_\ell) + (\mathbf{I}_\ell \otimes \mathbf{Q}).$$

According to the definition of \mathcal{N}_Q , every element $\text{vec}(\mathbf{Z})$ in \mathcal{N}_Q satisfies the condition

$$[\mathbf{K}_{\ell, m}(\mathbf{Q} \otimes \mathbf{I}_\ell) + (\mathbf{I}_\ell \otimes \mathbf{Q})]^\top \text{vec}(\mathbf{Z}) = \mathbf{0}$$

which leads to

$$\mathbf{Q}^\top \mathbf{Z} + \mathbf{Z}^\top \mathbf{Q} = \mathbf{0}.$$

This computation shows that $\mathcal{N}_Q \subset \{\text{vec}(\mathbf{Z}) : \mathbf{Z} \in \mathbf{R}^{m \times \ell}, \mathbf{Q}^\top \mathbf{Z} + \mathbf{Z}^\top \mathbf{Q} = \mathbf{0}\}$. The equality can be proved by computing the dimension of each space.

Since the columns of \mathbf{Q} and \mathbf{Q}_\perp can span \mathbf{R}^m , any matrix $\mathbf{Z} \in \mathbf{R}^{m \times \ell}$ can be represented as $\mathbf{Z} = \mathbf{Q}\mathbf{B} + \mathbf{Q}_\perp \mathbf{C}$, where $\mathbf{B} \in \mathbf{R}^{\ell \times \ell}$ and $\mathbf{C} \in \mathbf{R}^{(m-\ell) \times \ell}$. Another equivalent condition for null space is obtained by plugging this representation into the condition

$$\mathbf{Q}^\top (\mathbf{Q}\mathbf{B} + \mathbf{Q}_\perp \mathbf{C}) + (\mathbf{Q}\mathbf{B} + \mathbf{Q}_\perp \mathbf{C})^\top \mathbf{Q} = \mathbf{B} + \mathbf{B}^\top = \mathbf{0}.$$

Hence

$$\mathcal{N}_Q = \{\text{vec}(\mathbf{Q}\mathbf{B} + \mathbf{Q}_\perp \mathbf{C}) : \mathbf{B} \in \mathbf{R}^{\ell \times \ell}, \mathbf{C} \in \mathbf{R}^{(m-\ell) \times \ell}, \mathbf{B} = -\mathbf{B}^\top\}.$$

□

The second order necessary condition can be written out by the explicit form of the null space.

Lemma 3.3.4 (Second Order Necessary Condition). *Suppose \mathbf{Q} is a local maximizer of*

(+). Then \mathbf{Q} satisfies the inequality

$$\text{tr}(\mathbf{Z}^\top \bar{\mathbf{P}} \mathbf{Z} - \mathbf{Z}^\top \mathbf{Z} \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \leq 0.$$

for all $\text{vec}(\mathbf{Z}) \in \mathcal{N}_{\mathbf{Q}}$.

Proof. The second order necessary condition is

$$\text{vec}(\mathbf{Z})^\top \nabla_{\mathbf{Q}\mathbf{Q}}^2 \mathcal{L}(\mathbf{Q}, \Lambda) \text{vec}(\mathbf{Z}) \leq 0$$

for all $\text{vec}(\mathbf{Z}) \in \mathcal{N}_{\mathbf{Q}}$. By the relation of Lagrange multiplier at the feasible point (3.3) as $\Lambda + \Lambda^\top = \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}$, the gradient of \mathcal{L} by \mathbf{Q} can be written as

$$\begin{aligned} \nabla_{\mathbf{Q}} \mathcal{L}(\mathbf{Q}, \Lambda) &= (\mathbf{I} \otimes \bar{\mathbf{P}}) \text{vec}(\mathbf{Q}) - [\mathbf{K}_{\ell, m}(\mathbf{Q} \otimes \mathbf{I}_\ell) + (\mathbf{I}_\ell \otimes \mathbf{Q})] \text{vec}(\Lambda) \\ &= (\mathbf{I} \otimes \bar{\mathbf{P}}) \text{vec}(\mathbf{Q}) - \text{vec}(\mathbf{Q} \Lambda^\top + \mathbf{Q} \Lambda) \\ &= (\mathbf{I} \otimes \bar{\mathbf{P}}) \text{vec}(\mathbf{Q}) - [(\Lambda^\top + \Lambda) \otimes \mathbf{I}_m] \text{vec}(\mathbf{Q}) \\ &= (\mathbf{I} \otimes \bar{\mathbf{P}}) \text{vec}(\mathbf{Q}) - (\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q} \otimes \mathbf{I}_m) \text{vec}(\mathbf{Q}) \end{aligned}$$

Hence the Hessian matrix is

$$\nabla_{\mathbf{Q}\mathbf{Q}}^2 \mathcal{L}(\mathbf{Q}, \Lambda) = \mathbf{I} \otimes \bar{\mathbf{P}} - \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q} \otimes \mathbf{I}_m$$

The second order necessary condition can be rewritten as

$$\begin{aligned} \text{vec}(\mathbf{Z})^\top \nabla_{\mathbf{Q}\mathbf{Q}}^2 \mathcal{L}(\mathbf{Q}, \Lambda) \text{vec}(\mathbf{Z}) &= \text{vec}(\mathbf{Z})^\top \text{vec}(\bar{\mathbf{P}} \mathbf{Z} - \mathbf{Z} \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \\ &= \text{tr}(\mathbf{Z}^\top \bar{\mathbf{P}} \mathbf{Z} - \mathbf{Z}^\top \mathbf{Z} \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \leq 0. \end{aligned}$$

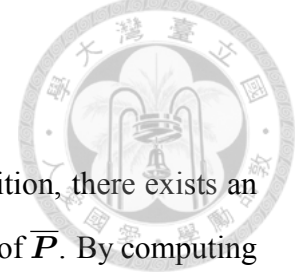
□

Finally, the main theorem of this part can be proved by the first and second order necessary conditions.

Theorem 3.3.5 (Local Maximizer of Optimization Problem). *Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_\ell > \lambda_{\ell+1} \geq \dots \geq \lambda_m$ be the eigenvalues of $\bar{\mathbf{P}}$ and let \mathbf{Q}_* be the corresponding leading ℓ*



eigenvectors. Then \mathbf{Q} is an local maximizer of (+) if and only if $\mathbf{Q} = \mathbf{Q}_* \mathbf{R}_*$ for some $\ell \times \ell$ orthogonal matrix \mathbf{R}_* .



Proof. By Lemma 3.3.2, for all \mathbf{Q} that satisfies the first order condition, there exists an orthogonal matrix \mathbf{R} such that each column of \mathbf{QR} is an eigenvector of $\bar{\mathbf{P}}$. By computing eigenvalue value decomposition,

$$\bar{\mathbf{P}} = (\mathbf{QR})\mathbf{S}(\mathbf{QR})^\top + \mathbf{Q}_\perp \mathbf{S}_\perp \mathbf{Q}_\perp^\top$$

where \mathbf{Q}_\perp contains $(m - \ell)$ eigenvector that orthogonal to \mathbf{QR} , $\mathbf{S} = \text{diag}(s_1, s_2, \dots, s_\ell)$, $\mathbf{S}_\perp = \text{diag}(s_1^\perp, s_2^\perp, \dots, s_{m-\ell}^\perp)$ and $s_1, s_2, \dots, s_\ell, s_1^\perp, s_2^\perp, \dots, s_{m-\ell}^\perp$ form all eigenvalues of $\bar{\mathbf{P}}$. By Lemma 3.3.3, the element $(\vec{\mathbf{Z}}) \in \mathcal{N}_{\mathbf{Q}}$ can be represented as $\mathbf{Z} = \mathbf{QB} + \mathbf{Q}_\perp \mathbf{C}$ with $\mathbf{B} = -\mathbf{B}$. These relations lead to a representation of second order condition as

$$\begin{aligned} & \text{tr}(\mathbf{Z}^\top \bar{\mathbf{P}} \mathbf{Z} - \mathbf{Z}^\top \mathbf{Z} \mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q}) \\ &= \text{tr}((\mathbf{QB} + \mathbf{Q}_\perp \mathbf{C})^\top ((\mathbf{QR})\mathbf{S}(\mathbf{QR})^\top + \mathbf{Q}_\perp \mathbf{S}_\perp \mathbf{Q}_\perp^\top) (\mathbf{QB} + \mathbf{Q}_\perp \mathbf{C}) \\ & \quad - (\mathbf{QB} + \mathbf{Q}_\perp \mathbf{C})^\top (\mathbf{QB} + \mathbf{Q}_\perp \mathbf{C}) \mathbf{Q}^\top ((\mathbf{QR})\mathbf{S}(\mathbf{QR})^\top + \mathbf{Q}_\perp \mathbf{S}_\perp \mathbf{Q}_\perp^\top) \mathbf{Q}) \\ &= \text{tr}(\mathbf{B}^\top \mathbf{R} \mathbf{S} \mathbf{R}^\top \mathbf{B}) + \text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) - \text{tr}(\mathbf{B}^\top \mathbf{B} \mathbf{R} \mathbf{S} \mathbf{R}^\top) - \text{tr}(\mathbf{C}^\top \mathbf{C} \mathbf{R} \mathbf{S} \mathbf{R}^\top) \leq 0 \end{aligned}$$

By the properties of trace and the relation $\mathbf{B} = -\mathbf{B}^\top$, $\text{tr}(\mathbf{B}^\top \mathbf{B} \mathbf{R} \mathbf{S} \mathbf{R}^\top) = \text{tr}(\mathbf{B} \mathbf{R} \mathbf{S} \mathbf{R}^\top \mathbf{B}^\top) = \text{tr}(\mathbf{B}^\top \mathbf{R} \mathbf{S} \mathbf{R}^\top \mathbf{B})$. Also, since $\mathbf{R} \mathbf{R}^\top = \mathbf{I}_\ell$, $\text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) = \text{tr}(\mathbf{R} \mathbf{R}^\top \mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) = \text{tr}((\mathbf{C} \mathbf{R})^\top \mathbf{S}_\perp (\mathbf{C} \mathbf{R}))$

Hence the second order condition in Lemma 3.3.4 can be written as

$$\begin{aligned} & \text{tr}(\mathbf{B}^\top \mathbf{R} \mathbf{S} \mathbf{R}^\top \mathbf{B}) + \text{tr}(\mathbf{C}^\top \mathbf{S}_\perp \mathbf{C}) - \text{tr}(\mathbf{B}^\top \mathbf{B} \mathbf{R} \mathbf{S} \mathbf{R}^\top) - \text{tr}(\mathbf{C}^\top \mathbf{C} \mathbf{R} \mathbf{S} \mathbf{R}^\top) \\ &= \text{tr}((\mathbf{C} \mathbf{R})^\top \mathbf{S}_\perp (\mathbf{C} \mathbf{R})) - \text{tr}((\mathbf{C} \mathbf{R})^\top (\mathbf{C} \mathbf{R}) \mathbf{S}_\perp) \\ &= \sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_i^\perp c_{ij}^2 - \sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_j c_{ij}^2 \leq 0 \end{aligned}$$

for all $\mathbf{C} \mathbf{R} = [c_{ij}] \in \mathbf{R}^{m \times (m-\ell)}$.

Let \mathbf{Q}_* denotes the leading ℓ eigenvectors of $\bar{\mathbf{P}}$. Suppose \mathbf{QR} contains eigenvector other than leading ℓ singular vectors of $\bar{\mathbf{P}}$. Then $s_a^\perp > s_b$ for some a and b . Now pick $[c_{ij}]$

as $c_{ab} = 1$ and $c_{ij} = 0$ otherwise. Then

$$\sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_i^{\perp} c_{ij}^2 - \sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_j c_{ij}^2 = s_a^{\perp} - s_b > 0.$$



This inequality conflicts to the second order necessary condition. On the other hand, suppose $\mathbf{QR} = \mathbf{Q}_*$. Then $\mathbf{S} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{\ell})$ and $\mathbf{S}_{\perp} = \text{diag}(\lambda_{\ell+1}, \dots, \lambda_m)$. Hence for all $[c_{ij}] \in \mathbf{R}^{m \times (m-\ell)}$.

$$\sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_i^{\perp} c_{ij}^2 - \sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_j c_{ij}^2 \leq (\lambda_{\ell+1} - \lambda_{\ell}) \sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_j c_{ij}^2 \leq 0$$

by $(\lambda_{\ell+1} - \lambda_{\ell}) < 0$ and $\sum_{i=1}^{(m-\ell)} \sum_{j=1}^{\ell} s_j c_{ij}^2 \leq 0$. This inequality obeys the second order necessary condition.

To sum up, a feasible point \mathbf{Q} satisfies the first and second order necessary condition if and only if $\mathbf{Q} = \mathbf{Q}_* \mathbf{R}_*$ for some orthogonal matrix $\mathbf{R}_* = \mathbf{R}^{\top}$. Lemma 3.1.1 shows that $\mathbf{Q}_* \mathbf{R}_*$ is global maximizer, and hence local maximizer. \square

Theorem 3.3.5 provides that the only type of local maximizer is the orthogonal matrix that formed by the leading ℓ eigenvectors of the matrix $\overline{\mathbf{P}}$ (with some right orthogonal transformation), which is the integrated subspaces defined in (2.1) or (3.1). This result shows that if a line search can prevent the iterator from saddle points, it can find out the integrated subspace successfully.



Chapter 4

Integration Method

As shown in Chapter 2, the integrated subspace is the leading eigenvectors of the matrix \bar{P} , or the singular vectors of the matrix $\left[\mathbf{Q}_{[1]} | \mathbf{Q}_{[2]} | \cdots | \mathbf{Q}_{[N]} \right]$. Therefore, canonical SVD (for example, the SVD routine in MATLAB or LAPACK [2]) can be applied to solve this eigenvectors problem. However, the computational complexity of canonical SVD is $O(N^2 m \ell^2)$, which is much slower when N increasing. Hence in this chapter, we introduce some other methods to compute the integrated subspace with different computational complexity than $O(N^2 m \ell^2)$.

4.1 Line Search Type Method

Since the integrated subspace is defined by an optimization problem, it is reasonable to use the algorithm for solving the optimization problem for computing the integrated subspace. A classical method for solving the optimization problem is line search methods. Gradient descent method (in our case, gradient ascent method) is widely used among the line search methods. The update scheme for unconstraint gradient ascent method is

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \tau_k \nabla f(\mathbf{x}_k)$$

where f is the objective function and τ_k is the step size needed to search in the n -th iteration. In the viewpoint of line search, in each iteration, we define a curve $\gamma_k(\tau) =$

$\mathbf{x}_k + \tau_k \nabla f(\mathbf{x}_k)$. This curve is a straight line starting at the point \mathbf{x}_k and pointing in the direction $\nabla f(\mathbf{x}_k)$, which is the steepest ascent direction of the objective function at \mathbf{x}_k . The point for the next iteration is given by finding a suitable step size τ and defining $\mathbf{x}_{k+1} = \mathbf{x}_k + \tau \nabla f(\mathbf{x}_k)$. These explanations show the main concept of the gradient ascent method. It tries to find the optimal solution by walking along the curve with the steepest direction in each step with a suitable step size.

The main issue for this algorithm is to find a suitable step size τ . If the step size is too large, it may be difficult to reach the optimal solution because it may walk too far. However, if the step size is too small, it may need more iterations to reach the optimal solution. One of the popular methods for selecting step size is using the backtracking method with the Armijo-Wolfe condition (for maximization problem)

$$\begin{aligned} f(\mathbf{x}_k + \tau_k \mathbf{p}_k) &\geq f(\mathbf{x}_k) + \rho \tau_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ \mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \tau_k \mathbf{p}_k) &\leq \rho_2 \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \end{aligned}$$

where ρ, ρ_2 are parameters that can be chosen and $\mathbf{p}_k = \nabla f(\mathbf{x}_k)$ in gradient ascent. The second condition always holds due to the design of the algorithm. Hence we just need to focus on the first condition. We use the Armijo rule for simplicity in the following content. The backtracking method tries to find the step size τ_k that satisfies the Armijo rule by the following steps. First, assign an initial guess of step size, such as 1. Then test whether this step size satisfies the Armijo rule. Suppose not, multiply this step size by a positive constant $\beta < 1$ and test it again. Repeat these processes until the step size satisfies the Armijo rule and set it into τ_k .

For the constrained optimization problem, the gradient ascent method needs to be modified. Here we choose to use the viewpoint of the gradient ascent method on the Stiefel manifold. The Stiefel manifold is the manifold consisting of all $m \times \ell$ orthogonal matrices

$$\mathcal{S}_{m,\ell} = \{Q \in \mathbf{R}^{m \times \ell} : Q^\top Q = I\}$$

which is exactly our constraint in the optimization problem.

Define

$$F(\mathbf{Q}) = \frac{1}{2} \text{tr}(\mathbf{Q}^\top \bar{\mathbf{P}} \mathbf{Q})$$

as our objective function. We denote the gradient of F at \mathbf{Q} as

$$\mathbf{G}_F(\mathbf{Q}) = \bar{\mathbf{P}} \mathbf{Q}$$



where the equality can be computed easily. On the manifold, the gradient in the Euclidean space may not represent the direction (tangent) at a point. Hence we need to project the gradient to the tangent space on the manifold for further derivation.

Lemma 4.1.1. *The projected gradient of F onto the tangent space $\mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$ of Stiefel manifold $\mathcal{S}_{m,\ell}$ is*

$$\mathbf{D}_F(\mathbf{Q}) = (\mathbf{I} - \mathbf{Q} \mathbf{Q}^\top) \bar{\mathbf{P}} \mathbf{Q}.$$

Proof. First we find a necessary and sufficient condition for \mathbf{X} being in $\mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$. For all $\mathbf{X} \in \mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$, find a path $\Gamma(t)$ in $\mathcal{S}_{m,\ell}$ with $\Gamma(0) = \mathbf{Q}$ and $\Gamma'(0) = \mathbf{X}$. From $\Gamma(t)^\top \Gamma(t) = \mathbf{I}$, differentiate each side by t and take $t = 0$, we have

$$\mathbf{X}^\top \mathbf{Q} + \mathbf{Q}^\top \mathbf{X} = \mathbf{0}, \quad (4.1)$$

which gives a necessary condition for $\mathbf{X} \in \mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$. There are $\ell(\ell + 1)/2$ conditions for \mathbf{X} in (4.1) and the dimension of $\mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$ is $m\ell - \ell(\ell + 1)/2$, which means (4.1) is also a sufficient condition for $\mathbf{X} \in \mathcal{T}_{\mathbf{Q}} \mathcal{S}_{m,\ell}$. By taking vec to each sides of (4.1), we get the equality

$$[(\mathbf{Q}^\top \otimes \mathbf{I}_\ell) \mathbf{K}_{m,\ell} + (\mathbf{I}_\ell \otimes \mathbf{Q}^\top)] \text{vec}(\mathbf{X}) = \mathbf{0}.$$

Define $\mathbf{T} = \mathbf{K}_{\ell,m}(\mathbf{Q} \otimes \mathbf{I}_\ell) + (\mathbf{I}_\ell \otimes \mathbf{Q})$ and get $\mathbf{T}^\top \text{vec}(\mathbf{X}) = \mathbf{0}$. This shows that the tangent space (after vectorizing each elements) is contained in the null space of \mathbf{T}^\top . One can compute the rank of \mathbf{T} and shows that the null space of \mathbf{T}^\top is actually the tangent space. Hence the projection matrix onto the tangent space is given by $(\mathbf{I} - \mathbf{P}_T)$, where $\mathbf{P}_T = \mathbf{T}(\mathbf{T}^\top \mathbf{T})^+ \mathbf{T}^\top$ and $(\mathbf{T}^\top \mathbf{T})^+$ denoted the Moore-Penrose pseudo-inverse. With \mathbf{P}_T ,

D_F can be given via $\text{vec}(D_F) = (\mathbf{I} - P_T) \text{vec}(G_F)$. With some calculation, we have $T = (\mathbf{I}_\ell \otimes Q)(\mathbf{I}_{\ell^2} + K_{\ell,\ell})$ and thus

$$\begin{aligned} T^\top T &= (\mathbf{I}_{\ell^2} + K_{\ell,\ell})^\top (\mathbf{I}_\ell \otimes Q)^\top (\mathbf{I}_\ell \otimes Q) (\mathbf{I}_{\ell^2} + K_{\ell,\ell}) \\ &= (\mathbf{I}_{\ell^2} + K_{\ell,\ell})(\mathbf{I}_{\ell^2} + K_{\ell,\ell}) = 2(\mathbf{I}_{\ell^2} + K_{\ell,\ell}). \end{aligned}$$



Then the projection matrix P_T can be calculated as:

$$\begin{aligned} P_T &= T(T^\top T)^+ T^\top \\ &= (\mathbf{I}_\ell \otimes Q)(\mathbf{I}_{\ell^2} + K_{\ell,\ell}) \frac{1}{2} (\mathbf{I}_{\ell^2} + K_{\ell,\ell})^+ (\mathbf{I}_{\ell^2} + K_{\ell,\ell})^\top (\mathbf{I}_\ell \otimes Q)^\top \\ &= \frac{1}{2} (\mathbf{I}_\ell \otimes Q)(\mathbf{I}_{\ell^2} + K_{\ell,\ell})(\mathbf{I}_\ell \otimes Q^\top) \frac{1}{2} (\mathbf{I}_\ell \otimes Q Q^\top) + \frac{1}{2} (Q^\top \otimes Q) K_{m,\ell}. \end{aligned}$$

Hence, by $\text{vec}(D_F) = (\mathbf{I} - P_T) \text{vec}(G_F)$,

$$\begin{aligned} \text{vec}(D_F) &= (\mathbf{I}_{\ell^2} - \frac{1}{2} (\mathbf{I}_\ell \otimes Q Q^\top) - \frac{1}{2} (Q^\top \otimes Q) K_{m,\ell}) \text{vec}(G_F) \\ &= \text{vec}(G_F) - \frac{1}{2} (\mathbf{I}_\ell \otimes Q Q^\top) \text{vec}(G_F) - \frac{1}{2} (Q^\top \otimes Q) K_{m,\ell} \text{vec}(G_F) \\ &= \text{vec}(G_F) - \frac{1}{2} \text{vec}(Q Q^\top G_F) - \frac{1}{2} \text{vec}(Q G_F^\top Q) \end{aligned}$$

and D_F can be written as

$$D_F = \left(\mathbf{I} - \frac{1}{2} Q Q^\top \right) G_F - \frac{1}{2} Q G_F^\top Q. \quad (4.2)$$

Since we have the property $Q^\top G_F(Q) = G_F(Q)^\top Q$ here, we can get $D_F(Q) = (\mathbf{I} - Q Q^\top) G_F(Q)$. This completes the proof. \square

Algorithm 3 is rewritten from [1] by the notation in this thesis and the projected gradient as above. Note that the function $\bar{F}(M)$ denote the function that orthogonalize M first (hence this point is in the Stiefel manifold) and then plug into F .

The basic concept of Algorithm 3 is same as the gradient ascent method described previously. For each step, we find the projected gradient, which is the steepest direction on the manifold. Then we walk along this direction with a suitable step size decide by backtracking and Armijo-Wolfe condition. The difference is that we need to retract back

Algorithm 3 Integration of $\{\mathbf{Q}_{[i]}\}_{i=1}^N$ based on Armijo line search.

Require: $\mathbf{Q}_{[1]}, \mathbf{Q}_{[2]}, \dots, \mathbf{Q}_{[N]}$ (subspace matrices), \mathbf{Q}_{ini} (initial guess), $\tau_0 > 0$ (initial step size), $\beta \in (0, 1)$ (scaling parameter for step size searching), $\rho \in (0, 1)$ (parameters for step size searching)

Ensure: Integrated subspace matrix $\overline{\mathbf{Q}}$ based on Armijo line search

- 1: Initialize the current iterate $\mathbf{Q}_c \leftarrow \mathbf{Q}_{\text{ini}}$
- 2: **while** (not convergent) **do**
- 3: Compute the gradient on manifold $\mathbf{X} = (\mathbf{I}_m - \mathbf{Q}_c \mathbf{Q}_c^\top) \overline{\mathbf{P}} \mathbf{Q}_c$
- 4: Find the smallest integer $j \geq 0$ such that the following inequality holds:

$$\overline{F}(\mathbf{Q}_c + \tau_0 \beta^j \mathbf{X}) \geq F(\mathbf{Q}_c) + \tau_0 \beta^j \rho \|\mathbf{X}\|_F^2$$

- 5: Orthogonalize $(\mathbf{Q}_c + \tau_0 \beta^j \mathbf{X})$ (for example, by QR-decomposition or polar decomposition) as \mathbf{Q}_+
 - 6: Assign $\mathbf{Q}_c \leftarrow \mathbf{Q}_+$
 - 7: **end while**
 - 8: Output $\overline{\mathbf{Q}} = \mathbf{Q}_c$
-

to the manifold before we calculate the objective function and before the next iteration. In this algorithm, it uses orthogonalization as the retraction.

In [1], the convergent theorem is also provided. Here we translate the theorem and write down the most important part as Theorem 4.1.2. Note that in the theorem, a critical point x_* (the points that make projected gradient is $\mathbf{0}$) is stable, if for any neighborhood \mathcal{U} around x_* , there exists another neighborhood \mathcal{V} around x_* , such that if the initial value start in \mathcal{V} , then after any finite steps, the result will be in \mathcal{U} . A critical point is asymptotically stable if it is stable, and the condition also holds as the number of steps tends to infinity. A critical point is unstable if it is not stable.

Theorem 4.1.2 (Convergence theory for Algorithm 3). *Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_\ell > \lambda_{\ell+1} \geq \dots \geq \lambda_m$ be the eigenvalues of $\overline{\mathbf{P}}$ and let \mathbf{Q}_* be the corresponding leading ℓ eigenvectors. Algorithm 3 converge to the orthogonal matrix \mathbf{Q} with each column are eigenvectors of $\overline{\mathbf{P}}$. Among all the critical points, \mathbf{Q}_* (up to a right orthogonal transform) is the only asymptotically stable point with the linear convergent rate. Other critical points are unstable.*

This theorem gives the promise of the convergence. Also, our result for the property of optimization problem in Chapter 3 can support this result. The critical points are actually the points satisfying first order condition 3.2. Among these points, only \mathbf{Q}_* satisfies the

second order condition, which is related to the condition of stable points.

However, same as the traditional gradient ascent method, the convergent rate of the Algorithm 3 is linear. This could be a problem since it might take many steps to converge. Here we provide another algorithm. Algorithm 4 is rewritten by the notation and objective function in this thesis from the algorithm proposed by Wen and Yin [11].

This algorithm is based on gradient ascent method which is similar to Algorithm 3. However, some modification is introduced for this algorithm. The method to retract the point back to the manifold is changed in this modified algorithm. The searching path at current point \mathbf{Q}_c is defined as $\Gamma_{\mathbf{Q}_c}(\tau) = (\mathbf{I} - \frac{\tau}{2}\mathbf{M})^{-1}(\mathbf{I} + \frac{\tau}{2}\mathbf{M})\mathbf{Q}_c$, where $\mathbf{M} = \mathbf{G}\mathbf{Q}_c^\top - \mathbf{Q}_c\mathbf{G}^\top$ and $\mathbf{G} = \overline{\mathbf{P}}\mathbf{Q}_c$ is the gradient of the objective function in Euclidean space. By using Woodbury matrix identity, the searching path is same as

$$\Gamma_{\mathbf{Q}_c}(\tau) = \mathbf{Q}_c - \tau\mathbf{L}(\mathbf{I}_{2\ell} + \frac{1}{2}\tau\mathbf{R}^\top\mathbf{L})^{-1}\mathbf{R}^\top\mathbf{Q}_c \quad (4.3)$$

where $\mathbf{L} = [-\mathbf{G} \ \mathbf{Q}_c]$ and $\mathbf{R} = [\mathbf{Q}_c \ \mathbf{G}]$. This modification decreases the matrix size to compute inverse. As mentioned in [11], this path also satisfies some properties, hence Theorem 4.1.2 also holds for this algorithm.

Algorithm 4 also uses the Barzilai-Borwein step size [3] (BB step size) to accelerate this gradient method, and the nonmonotone strategy in [12] to prevent stuck in the local optimal points. The Armijo rule in the Algorithm 4 is

$$F(\Gamma_{\mathbf{Q}_c}(\tau)) \geq F(\Gamma_{\mathbf{Q}_c}(0)) + \tau\sigma \left. \frac{dF(\Gamma_{\mathbf{Q}_c}(t))}{dt} \right|_{t=0}.$$

The nonmonotone strategy in [12] modified the Armijo rule by

$$F(\Gamma_{\mathbf{Q}_c}(\tau)) \geq c + \tau\sigma \left. \frac{dF(\Gamma_{\mathbf{Q}_c}(\tau))}{dt} \right|_{t=0}$$

where c is a scalar updated in each iteration by $c \leftarrow (\eta\zeta c + F(\mathbf{Q}_+))/(\eta\zeta + 1)$ and then $\zeta \leftarrow \eta\zeta + 1$ for a given parameter η . The initial value of c is $F(\mathbf{Q}_{\text{ini}})$. This condition soften the Armijo rule and increase the chance to jump out the local minimum.

The BB step size is the simulation for Newton's method, in a more efficient way. They consider the update scheme $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{S}_k \mathbf{g}_k$. In Newton's method, the matrix \mathbf{S}_k is the Hessian matrix of objective function. They set the matrix $\mathbf{S}_k = \tau_k \mathbf{I}$, where τ_k is the step size needed to be computed. Also, they wish \mathbf{S}_k can approximately satisfy the secant condition $\Delta \mathbf{x} = \mathbf{S}_k \Delta \mathbf{g}$ or $\mathbf{S}_k \Delta \mathbf{x} = \Delta \mathbf{g}$ in quasi-Newton's method, where $\Delta \mathbf{x} = \mathbf{x}_k - \mathbf{x}_{k-1}$, $\Delta \mathbf{g} = \mathbf{g}_k - \mathbf{g}_{k-1}$. Hence the step size τ_k is determined by

$$\tau_k = \underset{\tau}{\operatorname{argmin}} \|\Delta \mathbf{x} - \tau \Delta \mathbf{g}\| \quad \text{or} \quad \tau_k = \underset{\tau}{\operatorname{argmin}} \|\tau \Delta \mathbf{x} - \Delta \mathbf{g}\|.$$

The solutions of these minimization problems are

$$\alpha_k = \frac{\langle \Delta \mathbf{x}, \Delta \mathbf{g} \rangle}{\langle \Delta \mathbf{g}, \Delta \mathbf{g} \rangle} \quad \text{or} \quad \alpha_k = \frac{\langle \Delta \mathbf{x}, \Delta \mathbf{x} \rangle}{\langle \Delta \mathbf{x}, \Delta \mathbf{g} \rangle}$$

These two step sizes are called BB step size. In our case, the BB step sizes are

$$\tau_{guess} = \frac{\operatorname{tr}(\mathbf{D}_1^\top \mathbf{D}_1)}{|\operatorname{tr}(\mathbf{D}_1^\top \mathbf{D}_2)|} \quad \text{or} \quad \frac{|\operatorname{tr}(\mathbf{D}_1^\top \mathbf{D}_2)|}{\operatorname{tr}(\mathbf{D}_2^\top \mathbf{D}_2)}$$

where $\mathbf{D}_1 = \mathbf{Q}_+ - \mathbf{Q}_c$ and $\mathbf{D}_2 = \mathbf{X}_+ - \mathbf{X}_c$. However, BB step size does not imply the convergence. Therefore, Algorithm 4 still needs to use back tracking method with Armijo-Wolfe conditions to ensure the convergence of the algorithm.

In each iteration, the most intensive calculation in Algorithm 4 is the multiplication $\mathbf{G} = \overline{\mathbf{P}} \mathbf{Q}_c$. This multiplication can be computed by $\sum_{i=1}^N \mathbf{Q}_{[i]}^\top \mathbf{Q}_c$ first and then $\mathbf{G} = \sum_{i=1}^N \mathbf{Q}_{[i]} \mathbf{Q}_{[i]}^\top \mathbf{Q}_c$. The first step contains N matrix multiplications with size $\ell \times m$ and $m \times \ell$. The second step contains N matrix multiplications with size $m \times \ell$ and $\ell \times \ell$. Hence the complexity for computing \mathbf{G} is $O(Nm\ell^2)$.

The remaining steps in each iteration can be computed by using \mathbf{G} . First, we consider the part for line search. The value of objective function can be computed by $F(\mathbf{Q}) = \operatorname{tr}(\mathbf{Q}^\top \mathbf{G})$, which is the summation of dot product of two $m \times \ell$ matrices. The complexity is $O(m\ell)$. The point on the curve $\Gamma_{\mathbf{Q}_c}(\tau)$ can be solved by a linear system with dimension ℓ matrix multiplication of size $m \times 2\ell$ and $2\ell \times m$, and $m \times 2\ell$ and $2\ell \times 2\ell$. The complexity



Algorithm 4 Integration of $\{\mathbf{Q}_{[i]}\}_{i=1}^N$ based on nonmonotone line search with BB step size.

Require: $\mathbf{Q}_{[1]}, \mathbf{Q}_{[2]}, \dots, \mathbf{Q}_{[N]}$ (subspace matrices), \mathbf{Q}_{ini} (initial guess), $\tau_0 > 0$ (initial step size), $\beta \in (0, 1)$ (scaling parameter for step size searching), $\rho \in (0, 1)$ (parameter for step size searching), $\eta \in (0, 1)$ (parameter for next step searching), τ_M, τ_m (maximum and minimum for predicting step size)

Ensure: Integrated subspace matrix $\overline{\mathbf{Q}}$ based on Armijo line search with BB step size

- 1: Initialize $\mathbf{Q}_c \leftarrow \mathbf{Q}_{\text{ini}}, \bar{\tau} \leftarrow \tau_0, \zeta = 1, c = F(\mathbf{Q}_c)$
- 2: **while** (not convergent) **do**
- 3: Compute the gradient in Euclidean space $\mathbf{G} = \overline{\mathbf{P}}\mathbf{Q}_c$
- 4: Set $\mathbf{L} = [-\mathbf{G} \ \mathbf{Q}_c]$ and $\mathbf{R} = [\mathbf{Q}_c \ \mathbf{G}]$.
- 5: Find the smallest integer $j \geq 0$ such that the following inequality holds:

$$F(\mathbf{\Gamma}_{\mathbf{Q}_c}(\bar{\tau}\beta^j)) \geq c + \bar{\tau}\beta^j\rho \|\mathbf{Q}_c\mathbf{G}^\top - \mathbf{G}\mathbf{Q}_c^\top\|_F^2$$

where $\mathbf{\Gamma}_{\mathbf{Q}_c}(\tau) = \mathbf{Q}_c - \tau\mathbf{L}(\mathbf{I}_{2\ell} + \frac{1}{2}\tau\mathbf{R}^\top\mathbf{L})^{-1}\mathbf{R}^\top\mathbf{Q}_c$

- 6: Assign $\mathbf{Q}_+ = \mathbf{Q}_d(\bar{\tau}\beta^j)$
- 7: Update $c \leftarrow (\eta\zeta c + F(\mathbf{Q}_+))/(\eta\zeta + 1)$ and then $\zeta \leftarrow \eta\zeta + 1$
- 8: Compute the differences $\mathbf{D}_1 = \mathbf{Q}_+ - \mathbf{Q}_c$ and $\mathbf{D}_2 = \mathbf{X}_+ - \mathbf{X}_c$, where

$$\begin{aligned} \mathbf{X}_c &= (\mathbf{I}_m - \mathbf{Q}_c\mathbf{Q}_c^\top)\overline{\mathbf{P}}\mathbf{Q}_c \\ \mathbf{X}_+ &= (\mathbf{I}_m - \mathbf{Q}_+\mathbf{Q}_+^\top)\overline{\mathbf{P}}\mathbf{Q}_+ \end{aligned}$$

- 9: Assign $\bar{\tau} \leftarrow \max(\min(\tau_{\text{guess}}, \tau_M), \tau_m)$, where

$$\tau_{\text{guess}} = \frac{\text{tr}(\mathbf{D}_1^\top\mathbf{D}_1)}{|\text{tr}(\mathbf{D}_1^\top\mathbf{D}_2)|} \text{ or } \frac{|\text{tr}(\mathbf{D}_1^\top\mathbf{D}_2)|}{\text{tr}(\mathbf{D}_2^\top\mathbf{D}_2)}$$

- 10: Assign $\mathbf{Q}_c \leftarrow \mathbf{Q}_+$
 - 11: **end while**
 - 12: Output $\overline{\mathbf{Q}} = \mathbf{Q}_c$
-

is $O(m\ell^2 + \ell^3)$. Hence the total complexity for line search is $O(I_{inner}(m\ell^2 + \ell^3))$, where I_{inner} denotes the number of iteration of line search (inner loop). Second, we consider the part for updating c . The main calculation of this part is the computation of the objective function. Hence the complexity of this part is $O(m\ell)$. Finally, we consider the part for computing BB step size. The main computation in this part is \mathbf{X}_+ , which need to calculate $\mathbf{G}_+ = \overline{\mathbf{P}}\mathbf{Q}_+$. However, this computation can be directly used in next iteration. Hence we compute the complexity of this part in the next iteration. For other components, they need matrix multiplication with complexity $O(m\ell^2)$. Also, it needs $O(m\ell)$ to compute each trace for computing τ_{guess} .

To sum up, suppose we use I_{WY} iteration to converge, then the computational complexity of Algorithm 4 is $I_{WY}(O(Nm\ell^2) + O(I_{inner}(m\ell^2 + \ell^3)) + O(m\ell) + O(m\ell^2 + m\ell))$, which is dominate by the term $O(I_{WY}Nm\ell^2)$, suppose I_{inner} is controlled. (This assumption is reasonable since I_{inner} is often restricted within some number.)

4.2 Kolmogorov-Nagumo-Type Average

Besides using the viewpoint from line search method, one can use another viewpoint from the average on Stiefel manifold to compute the integrated subspace. Algorithm 5 is inspired by the Kolmogorov-Nagumo-type average on Stiefel manifold [8, 5]. If a pair of retraction map and lifting map is defined, then the algorithm can be generate by a fixed point method scheme. In these paper, a retraction map at a point \mathbf{Q} on Stiefel manifold is defined as a map $\varphi_{\mathbf{Q}}^{-1} : \mathcal{T}_{\mathbf{Q}}\mathcal{S}_{m,\ell} \rightarrow \mathcal{S}_{m,\ell}$ from the tangent space at \mathbf{Q} to the Stiefel manifold and satisfies the following three conditions: (1) $\varphi_{\mathbf{Q}}^{-1}$ can be defined around $\mathbf{0} \in \mathcal{T}_{\mathbf{Q}}\mathcal{S}_{m,\ell}$, (2) $\varphi_{\mathbf{Q}}^{-1}(\mathbf{0}) = \mathbf{Q}$ and (3) $\left. \frac{d\varphi_{\mathbf{Q}}^{-1}(t\mathbf{X})}{dt} \right|_{t=0} = \mathbf{X}$. The corresponding lifting map is a map $\varphi_{\mathbf{Q}} : \mathcal{S}_{m,\ell} \rightarrow \mathcal{T}_{\mathbf{Q}}\mathcal{S}_{m,\ell}$ from the Stiefel manifold to the tangent space and satisfies $\varphi_{\mathbf{Q}}^{-1}(\varphi_{\mathbf{Q}}(\mathbf{W})) = \mathbf{W}$. The fixed point scheme defied by this pair of lifting and retraction map is

$$\mathbf{Q}_+ = \varphi_{\mathbf{Q}_c}^{-1} \left(\frac{1}{N} \sum_{i=1}^N \varphi_{\mathbf{Q}_c}(\mathbf{Q}_{[i]}) \right).$$

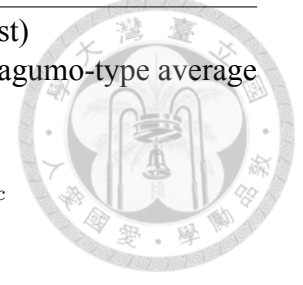
Note that \mathbf{Q}_+ is an Kolmogorove-Nagumo average of $\mathbf{Q}_{[i]}$.

Algorithm 5 Integration of $\{Q_{[i]}\}_{i=1}^N$ as Kolmogorov-Nagumo-type average.

Require: $Q_{[1]}, Q_{[2]}, \dots, Q_{[N]}$ (subspaces matrices), Q_{ini} (initial guest)

Ensure: Integrated subspace matrix \bar{Q} based on the Kolmogorov-Nagumo-type average

- 1: Initialize the current iterate $Q_c \leftarrow Q_{\text{ini}}$
 - 2: **while** (not convergent) **do**
 - 3: Perform the lifting map and average $\mathbf{X} = (\mathbf{I}_m - Q_c Q_c^\top) \bar{P} Q_c$
 - 4: Perform the retraction map $Q_+ \leftarrow Q_c C + \mathbf{X} C^{-1}$,
 where $C = \left\{ \frac{\mathbf{I}}{2} + \left(\frac{\mathbf{I}}{4} - \mathbf{X}^\top \mathbf{X} \right)^{1/2} \right\}^{1/2}$
 - 5: Assign $Q_c \leftarrow Q_+$
 - 6: **end while**
 - 7: Output $\bar{Q} = Q_c$
-



To connect the KN-type average to our integration problem (3.1), the lifting map is chosen as

$$\varphi_{Q_c}(\mathbf{W}) = (\mathbf{I}_m - Q_c Q_c^\top) \mathbf{W} \mathbf{W}^\top Q_c$$

which is an element in tangent space since $Q_c^\top \varphi_{Q_c}(\mathbf{W}) + \varphi_{Q_c}(\mathbf{W})^\top Q_c = \mathbf{0} + \mathbf{0} = \mathbf{0}$ for any $\mathbf{W} \in \mathcal{S}_{m,\ell}$. This choice makes the average of $Q_{[i]}$ on tangent space becomes

$$\sum_{i=1}^N \varphi_{Q_c}(Q_{[i]}) = \sum_{i=1}^N (\mathbf{I}_m - Q_c Q_c^\top) Q_{[i]} Q_{[i]}^\top Q_c = D_F(Q_c)$$

which is the projected gradient at Q_c . The problem now is what the corresponding retraction map is. As proposed in [4], the corresponding retraction map is chosen as

$$\varphi_{Q_c}^{-1}(\mathbf{X}) = Q_c C + \mathbf{X} C^{-1} \quad (4.4)$$

where $C = \left\{ \frac{\mathbf{I}}{2} + \left(\frac{\mathbf{I}}{4} - \mathbf{X}^\top \mathbf{X} \right)^{1/2} \right\}^{1/2}$. Algorithm 5 summarize these two maps with the fixed point scheme.

Here we give a quick derivation and check for the retraction map. For more detailed results, please refer to [4]. We want to find the retraction map $\varphi_{Q_c}^{-1}(\mathbf{X})$. Observe that from our choice of lifting map, we have the condition $Q_c^\top \varphi_{Q_c}(\mathbf{W}) = \mathbf{0}$ for any $\mathbf{W} \in \mathcal{S}_{m,\ell}$. Also, the condition we want for the retraction map is $\varphi_{Q_c}^{-1}(\varphi_{Q_c}(\mathbf{W})) = \mathbf{W}$. Hence the condition $\mathbf{X} = \varphi_{Q_c}(\mathbf{W})$ is a reasonable assumption. Also, assume $\varphi_{Q_c}^{-1}(\mathbf{X}) = Q_c C + \mathbf{X} B$. This assumption means that $\varphi_{Q_c}^{-1}(\mathbf{X})$ is spanned by the columns of Q_c and \mathbf{X} .

By using the condition $\varphi_{Q_c}^{-1}(\mathbf{X}) \in \mathcal{S}_{m,\ell}$ and $\varphi_Q^{-1}(\varphi_Q(\mathbf{W})) = \mathbf{W}$, we get the following equations

$$\begin{cases} \mathbf{I} = (\mathbf{Q}_c \mathbf{C} + \mathbf{X} \mathbf{B})^\top (\mathbf{Q}_c \mathbf{C} + \mathbf{X} \mathbf{B}) \\ \mathbf{W} = \mathbf{Q}_c \mathbf{C} + \mathbf{X} \mathbf{B} \\ \mathbf{X} = (\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \mathbf{W} \mathbf{W}^\top \mathbf{Q}_c \end{cases} .$$



Plug the second equation into third equation, we get $\mathbf{X} = \mathbf{X} \mathbf{B} \mathbf{C}^\top$. Here we add an assumption that \mathbf{C} is invertible and symmetric. Hence we have $\mathbf{X} \mathbf{B} = \mathbf{X} \mathbf{C}^{-1}$. Plug this back to the first equation and get $\mathbf{I} = \mathbf{C}^2 + \mathbf{C}^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{C}^{-1}$. This equation leads to $\mathbf{C}^2 = \mathbf{C}^4 + \mathbf{X}^\top \mathbf{X}$. Square this equation and get $(\mathbf{C}^2 - \frac{\mathbf{I}}{2})^2 = \frac{\mathbf{I}}{4} - \mathbf{X}^\top \mathbf{X}$, and hence $\mathbf{C} = \left\{ \frac{\mathbf{I}}{2} + \left(\frac{\mathbf{I}}{4} - \mathbf{X}^\top \mathbf{X} \right)^{1/2} \right\}^{1/2}$. To sum up, our guessed retraction map is same as the one given in (4.4).

No matter how many assumptions that we use for guessing the retraction map, we just need to check whether (4.4) is a corresponding retraction. We need to check the three conditions for retraction first. The guessed retraction map $\varphi_{Q_c}^{-1}(\mathbf{X})$ can be defined around $\mathbf{0} \in \mathcal{T}_Q \mathcal{S}_{m,\ell}$ since \mathbf{C} can be defined for $\text{tr}(\mathbf{X}^\top \mathbf{X}) \leq \frac{1}{4}$. The guessed retraction map satisfies $\varphi_Q^{-1}(\mathbf{0}) = \mathbf{Q}$ by directly computation. To check the third condition, we write $\varphi_Q^{-1}(t\mathbf{X}) = \mathbf{Q} \mathbf{C}(t) + t \mathbf{X} \mathbf{C}^{-1}(t)$. Derivative t to the both sides of condition $\mathbf{C}(t)^4 - \mathbf{C}(t)^2 + t^2 \mathbf{X}^\top \mathbf{X} = \mathbf{0}$, we get

$$\begin{aligned} & \frac{d\mathbf{C}(t)}{dt} \mathbf{C}(t)^3 + \mathbf{C}(t) \frac{d\mathbf{C}(t)}{dt} \mathbf{C}(t)^2 + \mathbf{C}(t)^2 \frac{d\mathbf{C}(t)}{dt} \mathbf{C}(t) + \mathbf{C}(t)^3 \frac{d\mathbf{C}(t)}{dt} \\ & + \frac{d\mathbf{C}(t)}{dt} \mathbf{C}(t) + \mathbf{C}(t) \frac{d\mathbf{C}(t)}{dt} + 2t \mathbf{X}^\top \mathbf{X} = \mathbf{0}. \end{aligned}$$

When $t = 0$, $\mathbf{C}(0) = \mathbf{I}$ and hence $\left. \frac{d\mathbf{C}(t)}{dt} \right|_{t=0} = \mathbf{0}$. Hence we can get $\left. \frac{d\varphi_Q^{-1}(t\mathbf{X})}{dt} \right|_{t=0} = \mathbf{X}$.

Now we check the condition that $\varphi_Q^{-1}(\varphi_Q(\mathbf{W})) = \mathbf{W}$. Here we may assume $\mathbf{W}^\top \mathbf{Q}$ is symmetric. Suppose not, compute the SVD of $\mathbf{W}^\top \mathbf{Q} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$ and replace \mathbf{W} by $\mathbf{W} \mathbf{U} \mathbf{V}^\top$. Let $\mathbf{X} = \varphi_Q(\mathbf{W}) = (\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \mathbf{W} \mathbf{W}^\top \mathbf{Q}_c$. Hence we get the equation

$$\mathbf{X}^\top \mathbf{X} = (\mathbf{W}^\top \mathbf{Q})^2 - (\mathbf{W}^\top \mathbf{Q})^4.$$

Solve for $(\mathbf{W}^\top \mathbf{Q})$ and get $\mathbf{W}^\top \mathbf{Q} = \mathbf{C}$. Hence we can directly compute

$$\begin{aligned}\varphi_{\mathbf{Q}}^{-1}(\varphi_{\mathbf{Q}}(\mathbf{W})) &= \mathbf{Q}_c \mathbf{C} + \mathbf{X} \mathbf{C}^{-1} \\ &= \mathbf{Q}_c \mathbf{C} + (\mathbf{I} - \mathbf{Q} \mathbf{Q}^\top) \mathbf{W} \mathbf{W}^\top \mathbf{Q} \mathbf{C}^{-1} \\ &= \mathbf{Q}_c \mathbf{C} + \mathbf{W} \mathbf{C} \mathbf{C}^{-1} - \mathbf{C}^2 \mathbf{C}^{-1} = \mathbf{W}\end{aligned}$$



Note that in the definition of integrated subspace, it is defined with orthogonal invariance of $\mathbf{Q}_{[i]}$. Hence it is reasonable to assume \mathbf{W} can be transformed by $\mathbf{U} \mathbf{V}^\top$ and satisfy the condition that $\mathbf{W}^\top \mathbf{Q}$ is symmetric.

An important issue in Algorithm 5 is the well definite of the matrix \mathbf{C} . We need to check that $(\frac{\mathbf{I}}{4} - \mathbf{X}^\top \mathbf{X})$ is semi positive definite for $\mathbf{X} = (\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \bar{\mathbf{P}} \mathbf{Q}_c$, and hence the square root of matrix can be defined (by SVD). It is equivalent to check $\|\mathbf{X}^\top \mathbf{X}\|_2 \leq \frac{1}{4}$. It is also equivalent to check $\|\mathbf{X}\|_2 \leq \frac{1}{2}$ by the relation $\|\mathbf{X}^\top \mathbf{X}\|_2 = \|\mathbf{X}\|_2^2$. By triangle inequality, we have

$$\|\mathbf{X}\|_2 \leq \frac{1}{N} \sum_{i=1}^N (\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \mathbf{Q}_{[i]} \mathbf{Q}_{[i]}^\top \mathbf{Q}_c.$$

Now we check that for every $\mathbf{W} \in \mathcal{S}_{m,\ell}$, the inequality $\|\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top\|_2 \leq \frac{1}{2}$ holds. It is equivalent to prove $\|(\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \mathbf{W} \mathbf{W}^\top \mathbf{Q}_c\|_2 \leq \frac{1}{4}$. Suppose the SVD of $\mathbf{Q}_c^\top \mathbf{W} = \mathbf{U} \mathbf{S} \mathbf{V}$. Then by directly computation

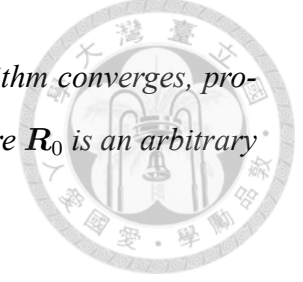
$$\begin{aligned}& \|(\mathbf{I} - \mathbf{Q}_c \mathbf{Q}_c^\top) \mathbf{W} \mathbf{W}^\top \mathbf{Q}_c\|_2 \\ &= \|\mathbf{U} \mathbf{S}^2 \mathbf{U} - \mathbf{U} \mathbf{S}^4 \mathbf{U}\|_2 \\ &= \|\mathbf{S}^2 - \mathbf{S}^4\|_2.\end{aligned}$$

Note that \mathbf{S} is diagonal matrix, and hence $(\mathbf{S}^2 - \mathbf{S}^4)$ is also diagonal matrix. For any real number x , we have the inequality $(x^2 - x^4) = -(x^2 - \frac{1}{2})^2 + \frac{1}{4} \leq \frac{1}{4}$. This equation shows that all diagonal entries of $(\mathbf{S}^2 - \mathbf{S}^4)$ is no greater than $\frac{1}{4}$. This result leads to $\|\mathbf{S}^2 - \mathbf{S}^4\|_2 \leq \frac{1}{4}$. The check is completed by tracing back the statements.

The convergence of Algorithm 5 is given as the following theorem. For the detailed

proof, please refer to [4].

Theorem 4.2.1. *There exists an $\varepsilon > 0$ such that KN-average algorithm converges, provided that the iteration starts from an initial $\mathbf{Q}_{\text{ini}} \in \mathcal{N}_\varepsilon(\mathbf{Q}_* \mathbf{R}_0)$, where \mathbf{R}_0 is an arbitrary orthogonal matrix.*



The Algorithm 5 can be also seen as a kind of gradient ascent method due to the selection of lifting map. Therefore, the computational complexity of this algorithm is similar to Algorithm 4. The dominant term of computational complexity is $O(I_{KN} N m \ell^2)$ if the algorithm needs I_{KN} iteration to converge. The dominant term is also comes from the computation of the projected gradient $\mathbf{X} = (\mathbf{I}_m - \mathbf{Q}_c \mathbf{Q}_c^\top) \overline{\mathbf{P}} \mathbf{Q}_c$, which is the average of lifting map for all sample subspaces $\mathbf{Q}_{[i]}$.

4.3 Reduction-Type Average

The main idea for this method is grouping the sample subspaces $\mathbf{Q}_{[i]}$ into several group, computing the integrated subspace of each groups, and then integrating these integrated subspaces. For example, suppose we have 4 sample subspaces need to be integrated and we choose the grouping number as 2. This method computes the leading ℓ left singular vectors of $[\mathbf{Q}_{[1]} \ \mathbf{Q}_{[2]}]$ as $\mathbf{Q}_{[1,2]}$, which is the integrated subspaces of $\mathbf{Q}_{[1]}$ and $\mathbf{Q}_{[2]}$. Next this method computes the leading ℓ left singular vectors of $[\mathbf{Q}_{[3]} \ \mathbf{Q}_{[4]}]$ as $\mathbf{Q}_{[3,4]}$. Then this method computes the average $\overline{\mathbf{Q}}$ as leading ℓ left singular vectors of $[\mathbf{Q}_{[1,2]} \ \mathbf{Q}_{[3,4]}]$. If the number of sample subspaces is more than 4, we can do it hierarchically.

The concept of reduction can be extended to any grouping number. However, we only consider the case that the grouping number is 2. The advantage of this case is that the leading left singular vectors of $\mathbf{M} = [\mathbf{Q}_1 \ \mathbf{Q}_2]$ can be written explicitly by the SVD of a small matrix $\mathbf{Q}_1^\top \mathbf{Q}_2 = \mathbf{U} \mathbf{S} \mathbf{V}^\top$. Suppose the SVD of $\mathbf{M} = \mathbf{L} \mathbf{\Sigma} \mathbf{R}^\top$ and the SVD of $\mathbf{Q}_1^\top \mathbf{Q}_2 = \mathbf{U} \mathbf{S} \mathbf{V}^\top$. The leading left singular vectors \mathbf{L}_ℓ can be obtained via the leading right singular vectors \mathbf{R}_ℓ and leading singular values Σ_ℓ as

$$\mathbf{L}_\ell = \mathbf{M} \mathbf{R}_\ell \Sigma_\ell.$$

Now, the problem is how to find the matrix \mathbf{R}_ℓ and Σ_ℓ . These two matrices can be obtained by the eigenvalue decomposition of $\mathbf{M}^\top \mathbf{M} = \mathbf{R}\Sigma^2\mathbf{R}^\top$. The singular value (eigenvalue) decomposition of $\mathbf{M}^\top \mathbf{M}$ can be write down explicitly by the SVD of $\mathbf{Q}_1^\top \mathbf{Q}_2 = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ as

$$\begin{aligned} \mathbf{M}^\top \mathbf{M} &= \begin{bmatrix} \mathbf{I}_\ell & \mathbf{Q}_1^\top \mathbf{Q}_2 \\ \mathbf{Q}_2^\top \mathbf{Q}_1 & \mathbf{I}_\ell \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I}_\ell & \mathbf{U}\mathbf{S}\mathbf{V}^\top \\ \mathbf{V}\mathbf{S}\mathbf{U}^\top & \mathbf{I}_\ell \end{bmatrix} \\ &= \left(\frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} & \mathbf{U} \\ \mathbf{V} & -\mathbf{V} \end{bmatrix} \right) \begin{bmatrix} \mathbf{I} + \mathbf{S} & \\ & \mathbf{I} - \mathbf{S} \end{bmatrix} \left(\frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} & \mathbf{U} \\ \mathbf{V} & -\mathbf{V} \end{bmatrix} \right)^\top. \end{aligned}$$

Hence $\Sigma_\ell = \mathbf{I} + \mathbf{S}$, $\mathbf{R}_\ell = \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}$ and

$$\mathbf{L}_\ell = \mathbf{M} \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix} (\mathbf{I} + \mathbf{S}) = (\mathbf{Q}_1\mathbf{U} + \mathbf{Q}_2\mathbf{V})(2(\mathbf{I} + \mathbf{S}))^{-\frac{1}{2}}.$$

Combine the above derivation and the hierarchical structure of the reduction, we can write down the algorithm as Algorithm 6.

Algorithm 6 Reduction

Require: The orthogonal matrices to be integrated $\mathbf{Q}_{[1]}, \mathbf{Q}_{[2]}, \dots, \mathbf{Q}_{[N]}$.

Ensure: The average $\overline{\mathbf{Q}}$.

- 1: Set $n = N$.
 - 2: **while** $n > 1$ **do**
 - 3: Set $m = \lfloor \frac{n}{2} \rfloor$
 - 4: **for** $i = 1, 2, \dots, m$ **do**
 - 5: Find SVD of $\mathbf{Q}_{[i]}^\top \mathbf{Q}_{[i+m]}$ as $\mathbf{U}\mathbf{S}\mathbf{V}^\top$.
 - 6: $\mathbf{Q}_{[i]} \leftarrow (\mathbf{Q}_{[i]}\mathbf{U} + \mathbf{Q}_{[i+m]}\mathbf{V})(2(\mathbf{I} + \mathbf{S}))^{-\frac{1}{2}}$.
 - 7: **end for**
 - 8: $n \leftarrow \lceil \frac{n}{2} \rceil$
 - 9: **end while**
 - 10: $\overline{\mathbf{Q}} = \mathbf{Q}_{[1]}$.
-

Since we have the explicit form of the integration of two subspaces, the computational cost of reduction becomes very cheap. The computation of integration of each pair of subspaces only needs matrix multiplication for $m \times \ell$ and $\ell \times \ell$ matrices. The computational complexity of this part is $O(m\ell^2)$. Also, only an $\ell \times \ell$ SVD needs to be computed for each pair. The computational complexity of this part is $O(\ell^3)$. For each pair, the computational complexity is $O(m\ell^2 + \ell^3)$, which is dominated by $O(m\ell^2)$. Due to the hierarchical structure of reduction, there are only $(N - 1)$ pairs need to be applied these processes. Therefore, the computational complexity of reduction is $O(Nm\ell^2)$, which is lower than the line search methods and canonical SVD.

The intuition of this algorithm is that if each group can get a good integrated subspace, the final integrated subspaces is also a good subspace. Since $\mathbf{Q}_{[i]}$ are sketched from a same matrix \mathbf{A} , the last statement should be correct in some sense. However, there is no related theory so far. In theory, the leading ℓ left singular vectors of $[\mathbf{Q}_{[1,2]} \ \mathbf{Q}_{[3,4]}]$ is in general not the integrated subspace defined in (2.1), which is the left singular vectors of $[\mathbf{Q}_{[1]} \ \mathbf{Q}_{[2]} \ \mathbf{Q}_{[3]} \ \mathbf{Q}_{[4]}]$. Although this fact, reduction still gives a roughly integrated subspace in our numerical experiment. Also, in the experiment, we try to use the result from reduction as the initial value of the line search algorithm.





Chapter 5

Comparison of rSVD and iSVD

Some comparison of rSVD and iSVD with same sketching number is shown in this chapter. In the beginning, we shall explain the meaning of rSVD and iSVD with same sketching number. Suppose in iSVD, we generate N sample subspace with rank ℓ . These numbers mean the sample subspace $\mathbf{Q}_{[i]}$ is an $m \times \ell$ orthogonal matrix for $i = 1, 2, \dots, N$. To generate these subspaces, we need to compute $\mathbf{Y}_{[i]} = \mathbf{A}\mathbf{\Omega}_{[i]}$ for all i . In the block matrix form

$$\left[\mathbf{Y}_{[1]} | \mathbf{Y}_{[2]} | \cdots | \mathbf{Y}_{[N]} \right] = \mathbf{A} \left[\mathbf{\Omega}_{[1]} | \mathbf{\Omega}_{[2]} | \cdots | \mathbf{\Omega}_{[N]} \right].$$

Now we change the viewpoint from iSVD to rSVD. Suppose the random matrix for sketching is $\mathbf{\Omega} = \left[\mathbf{\Omega}_{[1]} | \mathbf{\Omega}_{[2]} | \cdots | \mathbf{\Omega}_{[N]} \right]$ which is an $n \times N\ell$ matrix. Then the sketched matrix is $\mathbf{Y} = \mathbf{A}\mathbf{\Omega} = \left[\mathbf{Y}_{[1]} | \mathbf{Y}_{[2]} | \cdots | \mathbf{Y}_{[N]} \right]$ which is same as the case in iSVD. This derivation shows that using N sample subspaces with rank k in iSVD is equivalent to using $N\ell$ sketches in rSVD. Actually, both of them need to compute the same number of sketches. This is the meaning of rSVD and iSVD with same sketching number.

In common case, rSVD finds the orthogonal matrix of \mathbf{Y} as \mathbf{Q} and compute the SVD of $\mathbf{Q}\mathbf{Q}^\top \mathbf{A}$ as an approximation of the low-rank SVD. However, one can also find the leading singular vectors of \mathbf{Y} as \mathbf{Q} instead of using all the information from \mathbf{Y} . This can reduce the computational cost for computing the SVD of $\mathbf{Q}\mathbf{Q}^\top \mathbf{A}$. In the following discussion and the numerical experiment, this technique is included when we write the term ‘iSVD and rSVD with same sketching number.’

Now we can compare iSVD and rSVD with same sketching number. For rSVD, after sketching, we need to find the leading singular vectors of

$$\left[\mathbf{Y}_{[1]} | \mathbf{Y}_{[2]} | \cdots | \mathbf{Y}_{[N]} \right]$$

as \mathbf{Q} , and then computing the SVD of $\mathbf{Q}\mathbf{Q}^\top \mathbf{A}$. In iSVD, after sketching, we need to orthogonalize $\mathbf{Y}_{[i]}$ into $\mathbf{Q}_{[i]}$, and then find the integrated subspace $\overline{\mathbf{Q}}$. As mentioned in the previous chapter, we need to find the leading singular vectors of the matrix

$$\left[\mathbf{Q}_{[1]} | \mathbf{Q}_{[2]} | \cdots | \mathbf{Q}_{[N]} \right]$$

as $\overline{\mathbf{Q}}$, and then computing the SVD of $\overline{\mathbf{Q}}\overline{\mathbf{Q}}^\top \mathbf{A}$. Therefore, the difference between iSVD and rSVD is finding the different subspace to approximate the original matrix \mathbf{A} .

Intuitively, the approximation by rSVD may be accurate than iSVD since the block matrix in iSVD is the orthogonalization of each block in rSVD, which may lose the information of length. The numerical result shows that rSVD with same sketching number is slightly better than iSVD for accuracy. However, in theory, iSVD still get the exact low-rank approximation of \mathbf{A} if the number of sample subspace N goes large.

Due to the similarity between rSVD and iSVD, we can explain rSVD in the view-point of integration same as iSVD. Similar to the averaging concept as Theorem 3.2.1, the corresponding theorem for rSVD can be described as the following theorem.

Theorem 5.0.1. *Let the SVD of \mathbf{A} be $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. Let $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$, where $\mathbf{\Omega}$ is randomly generated by i.i.d. standard normal entries. Then*

$$E \left[\mathbf{Y}_{[i]} \mathbf{Y}_{[i]}^\top \right] = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^\top.$$



Proof. By direct computation,

$$\begin{aligned}
 E [\mathbf{Y}_{[i]} \mathbf{Y}_{[i]}^\top] &= E [\mathbf{A} \boldsymbol{\Omega}_{[i]} \boldsymbol{\Omega}_{[i]}^\top \mathbf{A}^\top] \\
 &= \mathbf{A} E [\boldsymbol{\Omega}_{[i]} \boldsymbol{\Omega}_{[i]}^\top] \mathbf{A}^\top \\
 &= \mathbf{A} \mathbf{I}_n \mathbf{A}^\top = \mathbf{U} \boldsymbol{\Sigma}^2 \mathbf{U}^\top
 \end{aligned}$$



where $E [\boldsymbol{\Omega}_{[i]} \boldsymbol{\Omega}_{[i]}^\top] = \mathbf{I}_n$ is derived directly from the condition that $\boldsymbol{\Omega}$ is randomly generated by i.i.d. standard normal entries. \square

The difference between the theorem for iSVD and rSVD is the spectrum matrix of the expected value. For rSVD, the spectrum matrix is $\boldsymbol{\Sigma}^2$, which is directly related to the spectrum of \mathbf{A} . For iSVD, the spectrum matrix is $\boldsymbol{\Lambda}$, which we can just describe some properties in the theorem.

For computing the leading singular vectors in rSVD, one can also use Algorithm 4 (WY) by replacing the matrix $\overline{\mathbf{P}}$ as

$$\overline{\mathbf{P}} = \frac{1}{N} \sum_{i=1}^N \mathbf{Y}_{[i]} \mathbf{Y}_{[i]}^\top.$$

However, Algorithm 5 (KN) can not be applied directly since the well definite of the retraction map will use the orthogonal property of $\mathbf{Q}_{[i]}$. The idea of reduction can be applied to rSVD. However, there is no explicit form for two subspace $\mathbf{Y}_{[1]}$, $\mathbf{Y}_{[2]}$. It can not generate the fast algorithm as Algorithm 6.





Chapter 6

Numerical Experiment

In this section, the performance of iSVD will be tested through the numerical experiment. The codes for testing are implemented in MATLAB. The desired rank in all the test is $k = 10$ and the exact sampling rank is $\ell = 22$. If a test only needs to observe the performance without timing, it is run on the machines with the larger size of memory. If a test needs to record the timing result, it is run on the MacBook Pro (Mid. 2014). (Processor: 2.6 GHz Intel Core i5. 2 cores. 4 threads. Memory: 8 GB 1600 MHz DDR3). All the tests follows the same steps as Algorithm 1 (rSVD) or Algorithm 2 (iSVD), but may use the different integration method. Table 6.1 shows the abbreviation and detail information of the integration algorithm used in this section.

The test matrices used in this paper are modified from the test matrix in [10]. These matrices are generated by the form $\mathbf{A} = \mathbf{H}_m \mathbf{\Sigma} \mathbf{H}_n^\top$, where \mathbf{H}_m denotes the $m \times m$ Hadamard matrix (a kind of orthogonal matrix with all the entries are 1 or -1), \mathbf{H}_n denotes the $n \times n$ Hadamard matrix, $m = 2^d$, $n = 2m = 2^{d+1}$, and $\mathbf{\Sigma}$ is an $m \times n$ diagonal matrix. In this

svds	The MATLAB built-in command svds
WY	Algorithm 4 with parameter $\beta = 0.5, \rho = 10^{-4}, \eta = 0.85$. Convergent condition: $\ \mathbf{D}_F(\mathbf{Q}_c)\ _2 < 10^{-3}$ (iSVD) Convergent condition: $\ \mathbf{D}_F(\mathbf{Q}_c)\ _2 < 10^{-3} \text{tr}(\mathbf{Y}_{[1]}^\top \mathbf{Y}_{[1]})/\ell$ (rSVD) Initial value $\mathbf{Q}_{\text{ini}} = \mathbf{Q}_{[1]}$
KN	Algorithm 5. Convergent condition: $\ \mathbf{I} - \mathbf{C}\ _F < 10^{-5}$ Initial value $\mathbf{Q}_{\text{ini}} = \mathbf{Q}_{[1]}$
red.	Reduction, Algorithm 6

Table 6.1: Abbreviation and detail information of the algorithm used in this section.

thesis, Σ is given by setting its diagonal entries as

$$\sigma_{i,i} = \begin{cases} s^{\frac{i-1}{k}} & \text{if } i \leq k \\ \frac{s(m-i)}{m-k-1} & \text{otherwise} \end{cases}$$



for $s = 10^{-1}, 10^{-3}$ and $k = 10$. We use $\mathbf{A}_H(10^{-1})$ and $\mathbf{A}_H(10^{-3})$ to denote the matrix generated with the case $s = 10^{-1}$ and $s = 10^{-3}$ respectively.

The similarity between each column of approximate leading singular vectors $\hat{\mathbf{U}}$ and exact leading singular vectors \mathbf{U}_k are used to measure the accuracy of the approximate leading left singular vectors $\hat{\mathbf{U}}$ from each test with different methods. The arranged left singular vectors \mathbf{U} of the test matrix is exactly the Hadamard matrix \mathbf{H}_m due to the construction of test matrix. Therefore, we compute the similarity between i -th column by computing the absolute value of inner product between the i -th column of \mathbf{U}_k and the i -th column of $\hat{\mathbf{U}}$. The more inner product close to 1, the better the approximation.

6.1 Different Number of Sketched Subspaces

The purpose of this test is to compare the result for the different choice of the number of sketched subspaces N with the different test matrix $\mathbf{A}_H(10^{-1})$ and $\mathbf{A}_H(10^{-3})$.

Figure 6.1 shows the similarity for different N with test matrix $\mathbf{A}_H(10^{-1})$ and $\mathbf{A}_H(10^{-3})$. The accuracy increases as N increasing, which is coherent to the explanation of 3.2.1 in Chapter 3. For $\mathbf{A}_H(10^{-3})$, as N increasing, the boxes in box plot become shorter. This provides that the variance of the similarity due to the randomness in sketching step is reduced as the number of sampling subspaces goes large. For $\mathbf{A}_H(10^{-1})$, the decreasing of variance is not same as $\mathbf{A}_H(10^{-3})$. The reason may be the quality for single subspaces ($N = 1$) is not good enough. Hence it is hard to capture the singular vectors after the fifth one. However, the accuracy still improves as N increasing, which is still a good result and coherent to the explanation in Chapter 3.

Next, we study the relation of N and the convergent iteration number of WY. Figure 6.2 shows that no matter the size of the matrix, the number of iteration for convergence

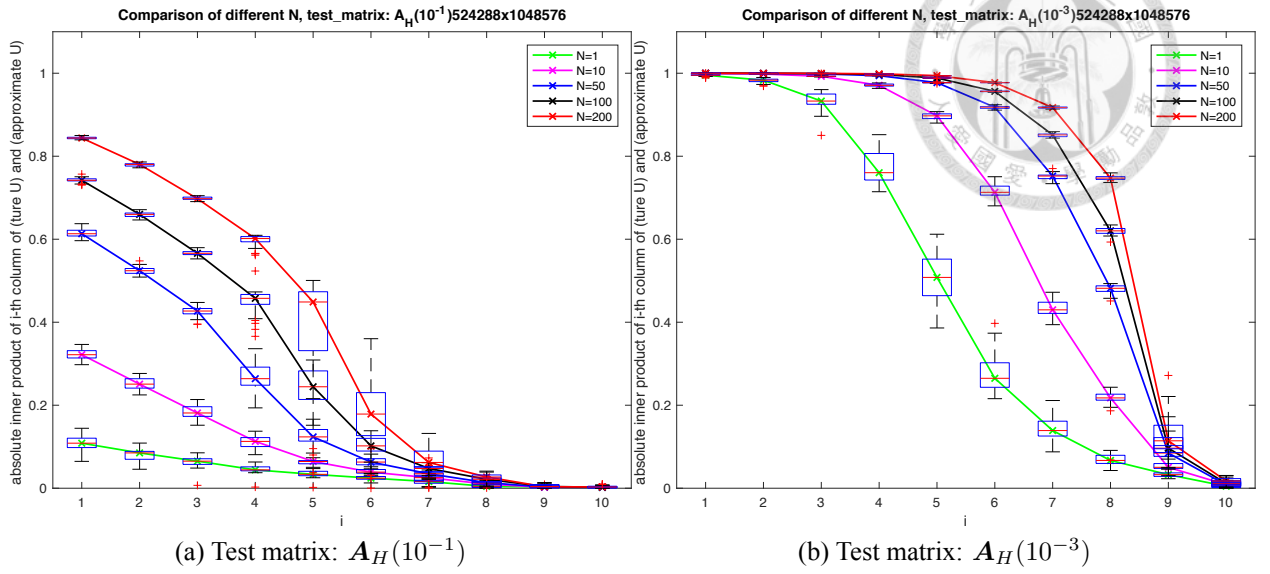


Figure 6.1: Similarity for different N . The size of test matrix is $m = 2^{19}$, $n = 2^{20}$. For each cases, we repeat 30 times iSVD with integration method WY and plot out the box plot of similarity. The box plot represent the maximum, Q3, median, Q1, minimum for each inner product among 30 times.

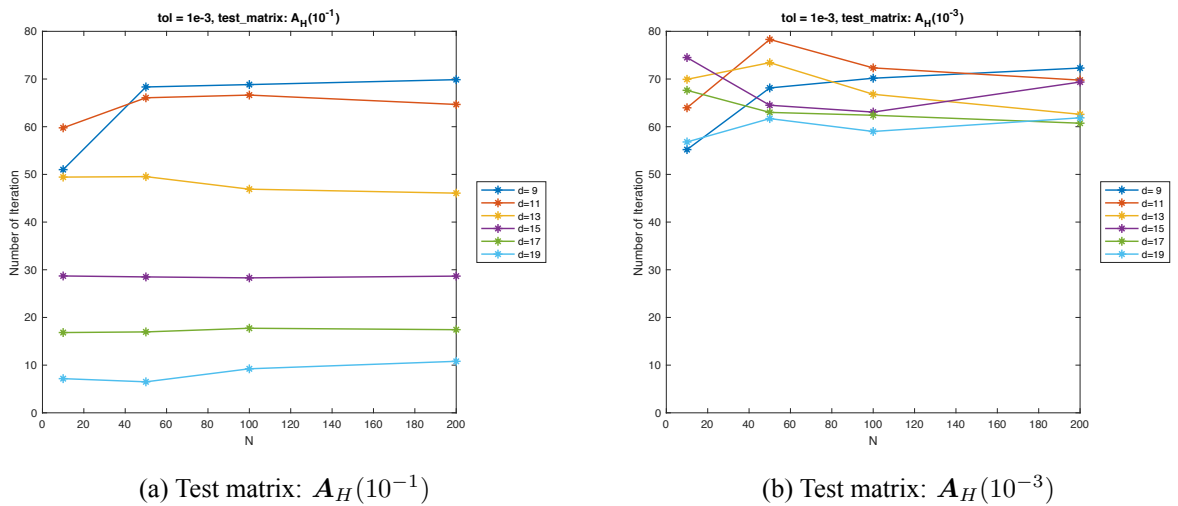


Figure 6.2: Average iteration number to converge for different N and different size of test matrix. The size of test matrix is $m = 2^d$, $n = 2^{d+1}$ for $d = 9, 11, 13, 15, 17, 19$. Each point shows the average iteration number among 30 tests.

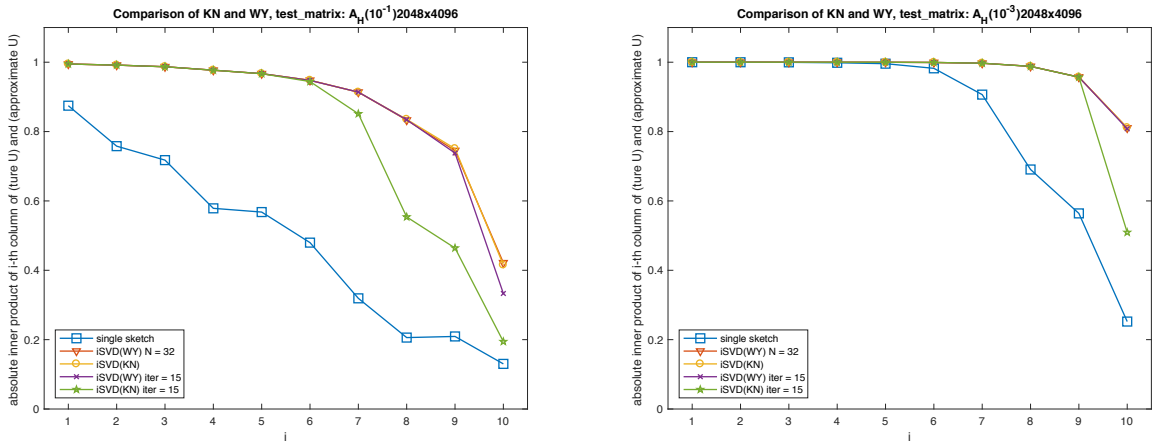
has no clear relation to the number of sampling subspaces N . This result is not surprising since the matrix $\overline{\mathbf{P}}$ tends to be a fixed matrix (expected value in Theorem 3.2.1) as N goes large. This is a good news since the computational complexity of WY is $O(I_{WY}Nm\ell^2)$ and the computational complexity of canonical SVD is $O(N^2m\ell^2)$. This figure shows I_{WY} do not increase as the number of sketched subspaces N increase. Hence using WY for integration is better in computational complexity than using canonical SVD.

6.2 Comparison of KN and WY

The purpose of this test is to compare the results of iSVD by using WY and KN. WY is derived from the viewpoint of line search and KN is derived from the viewpoint of average on Stiefel manifold. Although they are derived from different viewpoints, both of them contain the idea of gradient ascent. Therefore, it is interesting to observe the difference between these two algorithms.

Figure 6.3 shows the accuracy for the approximate singular vectors by using WY and KN. Both WY and KN can capture the approximate singular vectors with the same accuracy when they both converge. This result is coherent to the convergent theory in Section 4 that both of WY and KN converge to the same $\overline{\mathbf{Q}}$, and hence generate the same $\tilde{\mathbf{U}}$. However, KN needs more iteration to converge. To eliminate the effect of different convergence criteria, tow more results for the accuracy for WY and KN with same iteration number are added to the same figure. The approximation from WY is more accurate than KN when the iteration number is fixed as 15. This difference is not surprising since WY use BB step size to accelerate the convergence while KN does not use any technique to accelerate the convergence.

We do more test on the iteration number for WY and KN. Figure 6.4 and Figure 6.5 show the accuracy from WY and KN respectively, with the iteration number 5, 10, 15 and 20. As these tow figure shown, both WY and KN get a better accurate result as the iteration number goes large. WY converge faster than KN. For the test matrix $\mathbf{A}_H(10^{-1})$, WY can get almost same accurate as the converged result with only 20 iteration. For the test matrix $\mathbf{A}_H(10^{-3})$, WY only uses 15 iteration to get the almost same accurate result. In contrast,

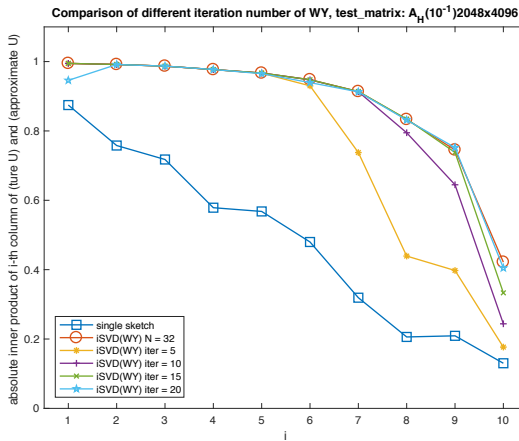


(a) Test matrix: $A_H(10^{-1})$

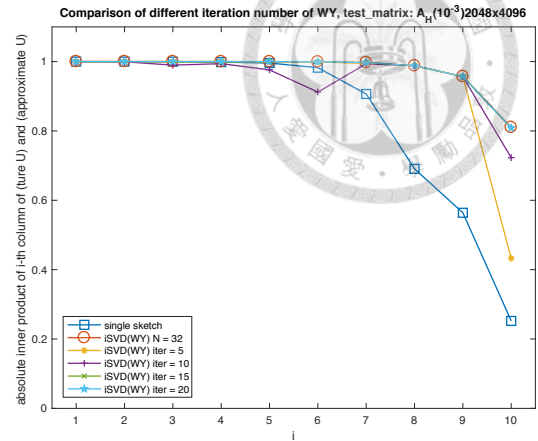
(b) Test matrix: $A_H(10^{-3})$

time(sec) iteration number	iSVD(WY)	iSVD(KN)	iSVD(WY, fix iter)	iSVD(KN, fix iter)
$A_H(10^{-1})$	0.964061 61	1.634885 242	0.266471 15	0.113637 15
$A_H(10^{-3})$	1.363975 84	1.538313 218	0.243105 15	0.142008 15

Figure 6.3: Comparison of the approximate singular vectors by using WY and KN. The size of test matrix is $m = 2^{11}$, $n = 2^{12}$ and the sampling number $N = 32$. All of these test use the same 32 sketched subspaces $Q_{[i]}$. The first line in the legend represents the similarity of $\bar{Q} = Q_{[1]}$. The second and third lines are the result from WY and KN respectively. The fourth and fifth lines are from WY and KN respectively with fixed iteration number 15.

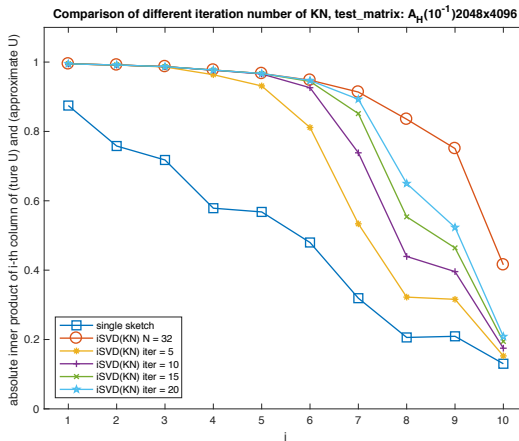


(a) Test matrix: $A_H(10^{-1})$

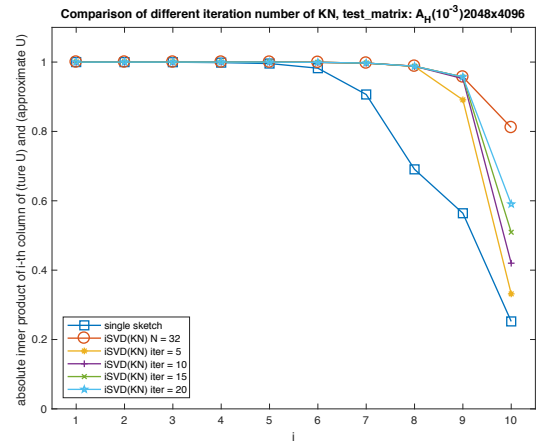


(b) Test matrix: $A_H(10^{-3})$

Figure 6.4: WY with different iteration numbers. The sketched subspaces $Q_{[i]}$ are same in Figure 6.3. The first line in the legend represents the similarity for the case $\bar{Q} = Q_{[1]}$. The second line is the similarity from WY (with 61 iteration for $A_H(10^{-1})$ and 84 iteration for $A_H(10^{-3})$ to converge). The third to sixth lines are from WY with iteration number 5, 10, 15, 20 respectively.



(a) Test matrix: $A_H(10^{-1})$



(b) Test matrix: $A_H(10^{-3})$

Figure 6.5: KN with different iteration numbers. The sketched subspaces $Q_{[i]}$ are same in Figure 6.3. The first line in the legend represents the similarity for the case $Q_{[1]}$. The second line is the similarity from KN. The third to sixth lines are from KN with iteration number 5, 10, 15, 20 respectively.

KN does not get the same accurate result in both cases by using 20 iteration.

However, the convergence of KN is more ‘smooth’ than WY. For the test matrix $\mathbf{A}_H(10^{-1})$, the similarity of the first column at 20 iteration is lower than the previous iteration. The same phenomenon also shows up in the sixth column of 10 iteration for the test matrix $\mathbf{A}_H(10^{-3})$. This could be a risk if one wants to use early stop technique in iSVD for WY. However, the fast convergence of WY could also be an advantage to use early stop technique.

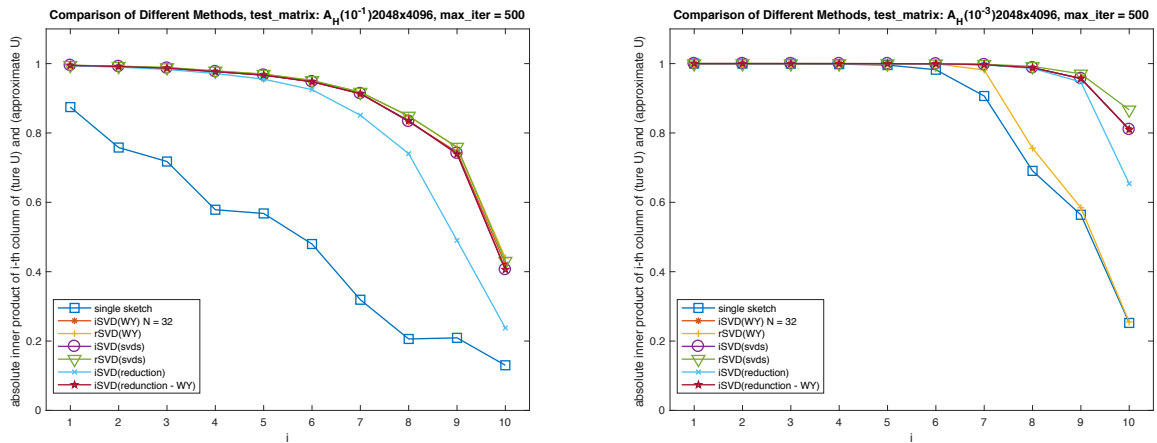
Also, these two figures point out that the stopping criterion for WY may not be suitable. For the test matrix $\mathbf{A}_H(10^{-3})$, WY can get almost same accurate as the converged result with only 20 iteration. However, the stopping criterion is not satisfied until 84 iteration. This phenomenon also happens in the test matrix $\mathbf{A}_H(10^{-1})$. The reason why this phenomenon happens may be the unsuitable choice of stopping criteria, since the stopping criteria measure whether \mathbf{Q}_c is convergent in WY, not measure whether \mathbf{Q}_c is an enough accurate subspace for generating the low-rank approximation of \mathbf{A} .

Remark. The timing results show that for a single iteration, WY uses about twice time than KN. However, the computation of complexity shows WY and KN have same computational complexity. This difference could be caused by the small size of the matrix and the lack of optimization for codes implemented WY.

6.3 Comparison of iSVD, rSVD and Reduction

The purpose of this test is to observe the difference between iSVD and rSVD numerically. As mentioned in the Chapter 5, iSVD and rSVD is very similar if the total sketching number is same. Also, we will do the numerical test of reduction in the same time to compare with rSVD more easily.

Figure 6.6 shows the result for iSVD and rSVD with WY and svds. We may treat the result from svds as the ideal result for iSVD and rSVD and compare it with the result from WY. For the test matrix $\mathbf{A}_H(10^{-1})$, the accuracy of rSVD is slightly better than iSVD. Also, both rSVD and iSVD with WY can capture the same accurate approximation as svds does. However, the iteration number of rSVD is more than the iteration number of

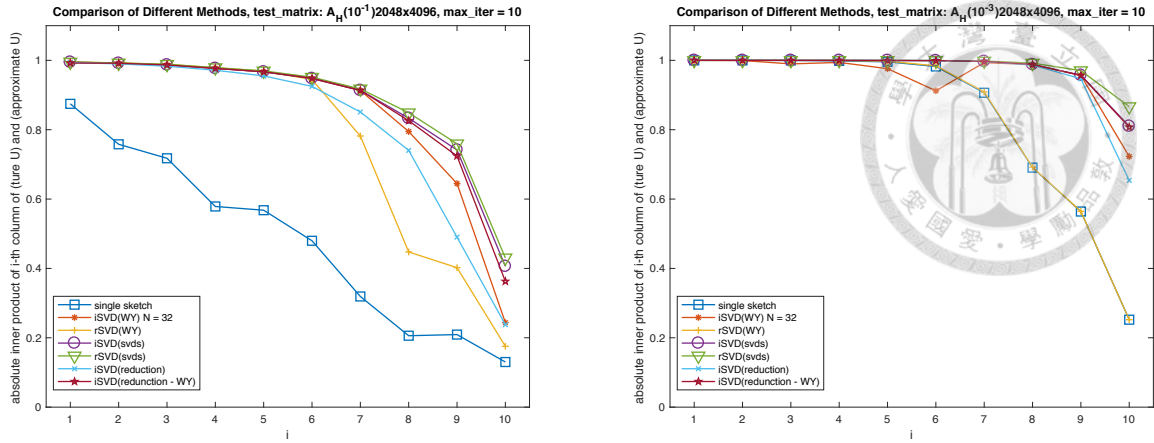


(a) Test matrix: $A_H(10^{-1})$

(b) Test matrix: $A_H(10^{-3})$

time(sec)	iSVD(WY)	rSVD(WY)	iSVD(svds)	rSVD(svds)	iSVD(red.)	iSVD (red.+WY)
iteration number						
$A_H(10^{-1})$	0.746769	1.839638	0.775641	0.692867	0.052223	0.901197
	61	122				66
$A_H(10^{-3})$	1.057405	0.865924	0.640089	0.597390	0.039458	1.017564
	84	63				79

Figure 6.6: Similarity for WY with iSVD and rSVD, reduction, and svds. The number of sketched subspaces in iSVD is $N = 32$. The number of sketching in rSVD is $32 * 22$, which is same as the total number of sketching in iSVD. The algorithm red.+WY uses the result of reduction as the initial value of WY.



(a) Test matrix: $A_H(10^{-1})$

(b) Test matrix: $A_H(10^{-3})$

time(sec) iteration number	iSVD(WY)	rSVD(WY)	iSVD(svds)	rSVD(svds)	iSVD(red.)	iSVD (red.+WY)
$A_H(10^{-1})$	0.158318 10	0.153800 10	0.823445	0.694192	0.039138	0.131893 10
$A_H(10^{-3})$	0.149794 10	0.166209 10	0.637057	0.617643	0.039255	0.133753 10

Figure 6.7: Similarity for the methods in 6.6 with the fixed iteration number 10 for WY.

iSVD. For the test matrix $A_H(10^{-3})$, rSVD with WY even fails to find the approximation with the same accuracy as svds. To eliminate the difference of convergent criteria, Figure 6.7 shows the results for the same setting with fixed iteration number 10. As shown in the figure, the accuracy for iSVD with WY is better than the accuracy for rSVD with WY when the iteration number is fixed. Although svds capture slightly better result than WY for iSVD, the computing time of WY is much faster than the svds.

For the reduction part, Figure 6.6 shows that reduction fails to capture the same accurate approximate as WY and svds. However, the computational time only takes about 0.04 second, which is very fast compare to WY and svds. Figure 6.7 shows that the accuracy for reduction is even better than the rSVD when the iteration number is 10.

Also, Figure 6.7 shows that using the result from reduction as the initial value of WY with only 10 iteration gives an almost accurate approximation same as svds does. However, the computational time for reduction+WY is much lower than svds. Although there is no theoretical guarantee for the error bound of the reduction so far, the numerical experiment shows that reduction can get an integrated subspace with better accuracy than

the single sketched subspaces, which is the original initial value for WY. It is not surprising that reduction+WY can get a better result than WY with same iteration number since reduction+WY start from a better initial guess.





Chapter 7

Discussion and Conclusion

In this thesis, some theoretical and performance analysis are shown. The integrated subspace \overline{Q} is defined by the solution of the optimization problem (2.1) or (3.1). The solution of these problems is the leading ℓ singular vectors of the matrix \overline{P} . Also, this is the only local maximizer of the problem (3.1).

Line search type algorithm (WY), KN type average (KN), and reduction are introduced to compute the integrated subspace. WY and KN are convergent to the integrated subspace. The leading term of computational complexity for WY and KN is $O(INm\ell^2)$, where I is the iteration number to converge. The computational complexity for canonical SVD for integration is $O(N^2m\ell^2)$. The numerical results show that N is independent of I . Also, WY can get nearly accurate approximation same as the convergent result by only few iteration number. These results provide that WY is an efficient algorithm for integration.

The computational complexity of reduction is $O(Nm\ell^2)$, which is faster than canonical SVD, WY, and KN. In theory, it does not capture the exact integrated subspace defined in (2.1). In the numerical experiment, the reduction can find an approximate integrated subspace. Also, it shows the potential that reduction can be used as a preprocess for WY to speed up the convergent with just a little extra computational cost.

If the total sketching number is same, iSVD is similar to rSVD. The difference between them is that the spectrum of the expected value in Theorem 3.2.1 and 5.0.1. Also, WY can directly apply to rSVD as the method to find the leading singular vectors of the sketching.

In the numerical experiment, WY converges faster for iSVD than rSVD. Sometimes WY also fails to find the leading singular vectors of the sketching in rSVD. However, there is no theoretical proof or explanation so far.

To sum up, iSVD gives an idea to integrate the subspaces generated from random sketching of a matrix. Although some phenomenons shown in numerical experiment lack theoretical explanation and proof, iSVD still shows the potential for approximate low-rank SVD with higher quality than rSVD in the same computing time. As more theoretical results showing up, iSVD could be an option for computing dimension reduction and feature extraction of large scale data faster but still accurate in the future.



Bibliography

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [2] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, third edition, 1999.
- [3] J. Barzilai and J. M. Borwein. Two-point step size gradient methods. *IMA Journal of Numerical Analysis*, 8(1):141–148, 1988.
- [4] T.-L. Chen, D. D. Chang, S.-Y. Huang, H. Chen, C. Lin, and W. Wang. Integrating multiple random sketches for singular value decomposition. *arXiv preprint arXiv:1608.08285*, 2016.
- [5] S. Fiori, T. Kaneko, and T. Tanaka. Mixed maps for learning a kolmogoroff-nagumo-type average element on the compact Stiefel manifold. *IEEE International Conference on Acoustic, Speech and Signal Processing (ICASSP)*, pages 4518–4522, 2014.
- [6] I. Griva, S. G. Nash, and A. Sofer. *Linear and nonlinear optimization*. Siam, 2009.
- [7] N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.
- [8] T. Kaneko, S. Fiori, and T. Tanaka. Empirical arithmetic averaging over the compact Stiefel manifold. *IEEE Transactions on Signal Processing*, 61(4):883–894, 2013.

- [9] J. R. Magnus and H. Neudecker. The commutation matrix: some properties and applications. *The Annals of Statistics*, pages 381–394, 1979.
- [10] V. Rokhlin, A. Szlam, and M. Tygert. A randomized algorithm for principal component analysis. *SIAM Journal on Matrix Analysis and Applications*, 31(3):1100–1124, 2009.
- [11] Z. Wen and W. Yin. A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 142(1-2):397–434, 2013.
- [12] H. Zhang and W. W. Hager. A nonmonotone line search technique and its application to unconstrained optimization. *SIAM Journal on Optimization*, 14(4):1043–1056, 2004.

