
New Probabilistic Bounds on Eigenvalues and Eigenvectors of Random Kernel Matrices

Nima Reyhani
Aalto University, School of Science
Helsinki, Finland
nima.reyhani@aalto.fi

Hideitsu Hino
Waseda University
Tokyo, Japan
hideitsu.hino@toki.waseda.jp

Ricardo Vigário
Aalto University, School of Science
Helsinki, Finland
ricardo.vigario@aalto.fi

Abstract

Kernel methods are successful approaches for different machine learning problems. This success is mainly rooted in using feature maps and kernel matrices. Some methods rely on the eigenvalues/eigenvectors of the kernel matrix, while for other methods the spectral information can be used to estimate the excess risk. An important question remains on how close the sample eigenvalues/eigenvectors are to the population values. In this paper, we improve earlier results on concentration bounds for eigenvalues of general kernel matrices. For distance and inner product kernel functions, e.g. radial basis functions, we provide new concentration bounds, which are characterized by the eigenvalues of the sample covariance matrix. Meanwhile, the obstacles for sharper bounds are accounted for and partially addressed. As a case study, we derive a concentration inequality for sample kernel target-alignment.

1 INTRODUCTION

Kernel methods such as Spectral Clustering, Kernel Principal Component Analysis (KPCA), and Support Vector Machines, are successful approaches in many practical machine learning and data analysis problems (Steinwart & Christmann, 2008). The main ingredient of these methods is the kernel matrix, which is built using the kernel function, evaluated at given sample points. The kernel matrix is a finite sample estimation of the kernel integral operator, which is determined by the kernel function. A number of algorithms rely on the eigenvalues and eigenvectors of the sample kernel operator/kernel matrix, and employ this information for further analysis. KPCA (Zwald & Blanchard, 2006) and kernel target-alignment (Cristianini

et. al., 2002) and (Jia & Liao, 2009) are good examples of such procedure. Also, the Rademacher complexity, and therefore the excess loss of the large margin classifiers can be computed using the eigenvalues of the kernel operator, see (Mendelson & Pajor, 2005) and references therein. Thus, it is important to know how reliably the spectrum of the kernel operator can be estimated, using a finite samples kernel matrix.

The importance of kernel and its spectral decomposition has initiated a series of studies of the concentration of the eigenvalues and eigenvectors of the kernel matrix around their expected values. Under certain rate of spectral decay of the kernel integral operator, (Koltchinskii, 1998) and (Koltchinskii & Giné, 2000) showed that the difference between population and sample eigenvalues and eigenvectors of a kernel matrix is asymptotically normal. A non-asymptotic analysis of kernel spectrum was first presented in (Shawe-Taylor et. al., 2005) and recently revisited in (Jia & Liao, 2009), where they derive an exponential bound for the concentration of the sample eigenvalues, using the bounded difference inequality. Also, using concentration results for isoperimetric random vectors, (Mendelson & Pajor, 2005) derived a concentration bound for the supremum for the concentration of eigenvalues. The bound is determined by the Orlicz norm of the kernel function with respect to the data distribution, the number of dimensions and the number of samples.

In this paper, we establish sharper exponential bounds for the concentration of sample eigenvalues of general kernel matrices, which depends on the minimum distance between a given eigenvalue and the rest of the spectrum. Separately, for Euclidean-distance and inner-product kernels, we provide a set of different bounds, which connect the concentration of eigenvalues to the maximum spectral gap of the sample covariance matrix. Similar results are derived for the eigenvectors of the same kernel matrix. Some experiments are also designed, to empirically test the presented re-

sults. As a case study, we derive concentration bounds for kernel target-alignment, which can be used to measure the agreement between a kernel matrix and the given labels.

The paper is organized as follows. Section 2 summarizes the previous results on concentration bounds for eigenvalues. In section 3, we present new concentration bounds for eigenvalues and eigenvectors of kernel matrices. Section 4 provides a concentration inequality for sample kernel alignment as a case study.

2 PREVIOUS RESULTS

In this section, we briefly present some of the most relevant earlier results on the concentration of the kernel matrix eigenvalues. The main assumption followed throughout the paper is summarized as

Assumption 1. Let $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^p$, be a set of samples of size n , independently drawn from distribution P . Let us define $K_n \in \mathbb{R}^{n \times n}$ to be a kernel matrix on \mathcal{S} , with $[K_n]_{i,j} = \frac{1}{n}k(\mathbf{x}_i, \mathbf{x}_j)$, for Mercer's kernel function k .

The kernel function k defines a kernel integral operator by

$$Tf(\cdot) = \int k(\mathbf{x}, \cdot)f(\mathbf{x})dP(\mathbf{x}),$$

for all smooth functions f . The eigenvalue equation is defined by $T\mathbf{u}(\cdot) = \lambda\mathbf{u}(\cdot)$, where λ is an eigenvalue of T and $\mathbf{u} : \mathbb{R}^p \rightarrow \mathbb{R}$ is the corresponding eigenfunction. For data samples \mathcal{S} , the kernel operator can be estimated by the kernel matrix K_n . Then, the eigenvalue problem can be reduced to

$$K_n\mathbf{u} = \lambda\mathbf{u}.$$

The solutions of the above equation are denoted by $\lambda_i(K_n)$ and $\mathbf{u}_i(K_n)$, $i = 1, \dots, n$. For simplicity of the paper, we assume the eigenvalues of kernel operator/matrix are not identical and are indexed in decreasing order, i.e. $\lambda_1 > \lambda_2 > \dots$. The focus of this paper is mainly in finding bounds for the concentrations of type: $P\{|\frac{1}{n}\lambda_i(K_n) - \frac{1}{n}\mathbb{E}_{\mathcal{S}}\lambda_i(K_n)| \leq \epsilon\}$ and similarly $P\{\|\frac{1}{n}\mathbf{u}_i(K_n) - \frac{1}{n}\mathbb{E}_{\mathcal{S}}\mathbf{u}_i(K_n)\| \leq \epsilon\}$, for all $i = 1, \dots, n$. We denote the finite dimensional unit sphere in \mathbb{R}^p by \mathcal{S}^{p-1} , i.e. $\mathcal{S}^{p-1} := \{\mathbf{w} : \mathbf{w} \in \mathbb{R}^p, \|\mathbf{w}\|_2 = 1\}$, and the Lipschitz norm of function f by $|f|_L$.

The following theorem is presented in (Shawe-Taylor et al., 2005) and provides a uniform bound that depends only on the supremum of the diagonal elements of the kernel matrix. This theorem requires that the kernel is bounded.

Theorem 2.1 ((Shawe-Taylor et al. 2005)). *Let us*

take Assumption 1. Then, for all $\epsilon > 0$, we have

$$P\{|\frac{1}{n}\lambda_i(K_n) - \mathbb{E}_{\mathcal{S}}\frac{1}{n}\lambda_i(K_n)| > \epsilon\} \leq 2\exp(-\frac{2n\epsilon^2}{R^4}),$$

where $R^2 = \max_{\mathbf{x} \in \mathbb{R}^p} k(\mathbf{x}, \mathbf{x})$.

The bound provided by Theorem 2.1 is suboptimal, as it only depends on the kernel function k through R and also it is uniform for all different eigenvalues. To address the suboptimality of Theorem 2.1, (Jia & Liao, 2009) proposed to bound the concentration of the kernel matrix eigenvalues using the largest eigenvalue of the sample kernel and a quantity $\theta \in (0, 1)$:

Theorem 2.2 ((Jia & Liao, 2009)). *Let us take Assumption 1. Then, for every $\epsilon > 0$, there exists $0 < \theta \leq 1$, such that*

$$P\{|\frac{1}{n}\lambda_i(K_n) - \mathbb{E}_{\mathcal{S}}\frac{1}{n}\lambda_i(K_n)| > \epsilon\} \leq 2\exp(-\frac{2\epsilon^2}{\theta^2\lambda_1^2(K_n)}).$$

The advantage of this result over the uniform result in Theorem 2.1 is that the concentration bound depends solely on a parameter θ and the largest sample eigenvalue of the kernel matrix. The parameter θ , in principle, varies as a function of the order of the eigenvalue. However, Theorem 2.2 does not provide any hint on estimating an optimal value for the parameter θ for any specific eigenvalue. In the following section, we provide sharper concentration inequalities for the spectral decomposition of the kernel matrix.

3 RESULTS

The main result is presented in Theorem 3.1, which establishes a bound on the concentration of eigenvalues of the kernel matrix, for any general kernel function that fulfills Mercer's theorem. In the next subsections, we provide concentration bounds for the spectral decomposition of the distance and inner product kernels. Our first result is inspired by the following observation: we built kernel matrices, using Gaussian kernels from 100 samples drawn from a 5-dimensional multivariate Gaussian. We repeat this example 1000 times. The boxplot of the first 15 eigenvalues of the Gaussian kernel matrix are depicted in Figure 1. From the boxplot, we can see that the distance between any box's corner and location of median decreases as the gap between corresponding eigenvalue and the spectrum of the kernel matrix decreases. This suggests that the concentration of the eigenvalues might be controlled by the distance of that eigenvalue to the spectrum. None of the earlier results, i.e. Theorems 2.1 and 2.2, provide such a characterization. Theorem 3.1 establishes a concentration bound, which illustrates the aforementioned phenomenon.

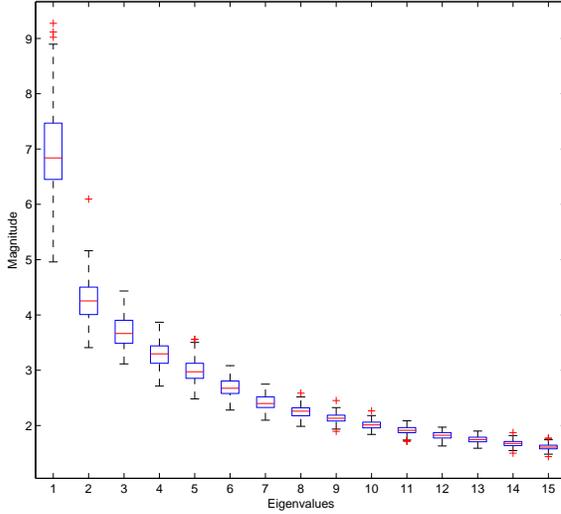


Figure 1: Boxplot of 15 first eigenvalues of Gaussian kernels with 5 dimensional Gaussian samples, drawn from $\mathcal{N}(0, I_5)$.

Theorem 3.1 (Probabilistic bound on eigenvalues of kernel matrices). *Let us take (Assumption 1) and define $\lambda_{i,i+1}(K_n) := \lambda_i(K_n) - \lambda_{i+1}(K_n)$. Then, the following inequality holds for sample eigenvalues, $\lambda_i(K_n)$, $i = 1, \dots, n - 1$.*

$$P\left\{\left|\frac{1}{n}\lambda_i(K_n) - \mathbb{E}_S\frac{1}{n}\lambda_i(K_n)\right| > \epsilon\right\} \leq \exp\left(-\frac{2n\epsilon^2}{\lambda_{i,i+1}^2(K_n)}\right). \quad (1)$$

Experimental results to validate the above bound are given in Example 1. To prove Theorem 3.1, we introduce Cauchy’s interlacing lemma, as a particular matrix perturbation result and bounded difference inequality.

Definition 1 (Principal Submatrix). *The principal submatrix of A of order $n - 1$ is the matrix obtained by taking the first $n - 1$ rows and columns of matrix A .*

Lemma (Cauchy’s Interlacing lemma, (Horn & Johnson, 1985)). *Let B be a principal submatrix of the Hermitian matrix A , of order $n - 1$, with eigenvalues $\mu_1 \geq \dots \geq \mu_{n-1}$. Also, we denote the eigenvalues of A by $\lambda_1 \geq \dots \geq \lambda_n$. Then, the following property holds for eigenvalues of A and B :*

$$\lambda_1 \geq \mu_1 \geq \lambda_2 \geq \mu_2 \geq \dots \geq \mu_{n-1} \geq \lambda_n.$$

Theorem (Bounded Difference Inequality, (Ledoux & Talagrand, 1991) and (Shawe-Taylor et. al., 2005)). *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be independent random variables taking values in a set \mathcal{A} and $g : \mathcal{A}^n \rightarrow \mathbb{R}$ satisfies one of the*

following: for $c_i \geq 0, 1 \leq i \leq n$,

$$\sup_{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}'_i} |g(\mathbf{x}_1, \dots, \mathbf{x}_n) - g(\mathbf{x}_1, \dots, \mathbf{x}'_i, \mathbf{x}_{i+1}, \mathbf{x}_n)| \leq c_i$$

$$\sup_{\mathbf{x}_1, \dots, \mathbf{x}_n} |g(\mathbf{x}_1, \dots, \mathbf{x}_n) - g_i(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \mathbf{x}_n)| \leq c_i,$$

where $\mathbf{x}'_1, \dots, \mathbf{x}'_n$ are independent copies of $\mathbf{x}_1, \dots, \mathbf{x}_n$ and $g_i : \mathcal{A}^{n-1} \rightarrow \mathbb{R}, i = 1, \dots, n$. Then, for all $\epsilon > 0$ we have $P\{|g(\mathbf{x}_1, \dots, \mathbf{x}_n) - \mathbb{E}g(\mathbf{x}_1, \dots, \mathbf{x}_n)| > \epsilon\} \leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}\right)$.

Proof of Theorem 3.1. Let us define submatrix $K_n^n \in \mathbb{R}^{n-1 \times n-1}$ with $[K_n^n]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j), i, j = 1, \dots, n - 1$. By the Cauchy’s interlacing lemma, we have $\lambda_i(K_n) - \lambda_i(K_n^n) \leq \lambda_i(K_n) - \lambda_{i+1}(K_n) = \lambda_{i,i+1}(K_n)$. Using the bounded difference inequality, see Appendix 1, and taking $\sum_{j=1}^n c_j^2 = n\lambda_{i,i+1}^2(K_n)$ we obtain the claimed result. \square

Note that, $\lambda_{i,i+1}(K_n)$ can be bounded by $\theta\lambda_1$ for some $\theta \in (0, 1]$, which reduces to the bound in Theorem 2.2. With extra conditions on the condition number of kernel matrix we can derive the bound in Theorem 2.1 by the above result, for details please see Corollary 2 in (Jia & Liao, 2009).

Corollary 3.1. *Using the assumptions of Theorem 3.1, the sum of the top and low eigenvalues are concentrated around their mean, as follows:*

$$\begin{aligned} P\left\{\left|\frac{1}{n}\sum_{i=1}^k \lambda_i(K_n) - \mathbb{E}_S\frac{1}{n}\sum_{i=1}^k \lambda_i(K_n)\right| > \epsilon\right\} \\ \leq \exp\left(-\frac{2n\epsilon^2}{\lambda_{1,k+1}^2(K_n)}\right), \\ P\left\{\left|\frac{1}{n}\sum_{i=k}^n \lambda_i(K_n) - \mathbb{E}\frac{1}{n}\sum_{i=k}^n \lambda_i(K_n)\right| > \epsilon\right\} \\ \leq \exp\left(-\frac{2n\epsilon^2}{\lambda_{k,n}^2(K)}\right), \end{aligned}$$

where $\lambda_{1,k+1}(K_n) = \lambda_1(K_n) - \lambda_{k+1}(K_n)$ and $\lambda_{k,n}(K_n) = \lambda_k(K_n) - \lambda_n(K_n)$.

The above inequality is useful to derive a probabilistic bound for the population risk of Kernel PCA, using the empirical risk. We leave further details for future work.

3.1 DISTANCE AND INNER PRODUCT KERNELS

In the following, we present results for distance and inner product kernel functions, which are commonly used in practical data analysis, such as Radial Basis Functions in SVM (Steinwart & Christmann, 2008), or Laplacian kernels in Spectral Clustering (Bengio et.

al., 2003). The distance kernel matrix and inner product kernel matrix are defined by $k(\mathbf{x}_i, \mathbf{x}_j) = f(\|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$ and $k(\mathbf{x}_i, \mathbf{x}_j) = f(\mathbf{x}_i^\top \mathbf{x}_j)$, where $f: \mathbb{R} \rightarrow \mathbb{R}$ is a smooth function, which has a Bochner's integral representation (Steinwart & Christmann, 2008). For the rest of the paper, except section 4, the spectral norm is denoted by $\|\cdot\|$. The results for concentration of eigenvalues of the distance kernel matrix is presented in the Theorem 3.2.

Theorem 3.2. *Let us take (Assumption 1), with $k(\mathbf{x}_i, \mathbf{x}_j) = f(\|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$. We assume $\mathbf{x}_1, \dots, \mathbf{x}_n$ have a nonsingular sample covariance matrix Σ . Also we assume, for r.v.s $\mathbf{y}_1, \dots, \mathbf{y}_n$ with identity covariance matrix and same distribution as $\mathbf{x}_1, \dots, \mathbf{x}_n$, $\|\mathbf{y}_i\|_2 \leq M$ and $\mathbf{x}_i = \Sigma^{\frac{1}{2}} \mathbf{y}_i, \forall i = 1, \dots, n$ holds. The first and last eigenvalues of Σ are denoted by $\lambda_1(\Sigma)$ and $\lambda_p(\Sigma)$, respectively. Then, for $\lambda_{1,p} = \lambda_1(\Sigma) - \lambda_p(\Sigma)$, we have*

$$P\left\{\left|\frac{1}{n}\lambda_i(K) - \mathbb{E}_S\frac{1}{n}\lambda_i(K)\right| > \epsilon\right\} \leq \exp\left(-\frac{n^2\epsilon^2}{18M^4|f|_L^2\lambda_{1,p}^2}\right). \quad (2)$$

Proof. Let us denote the perturbation of matrix K by K^n defined by $K_{i,j}^n := \frac{1}{n}f(\|\mathbf{x}_i - \mathbf{x}_j\|_2^2), \forall i, j = 2, \dots, n$ and $K_{j,1}^n = K_{1,j}^n := \frac{1}{n}f(\|\mathbf{x}_{1'} - \mathbf{x}_j\|), \forall j = 1', \dots, n$, where $\mathbf{x}_{1'}$ is independent copy of $\mathbf{x}_i, \forall i = 1, \dots, n$. Then, $K = K^n + E$. Suppose the perturbation term E has a small spectral norm. Then, by the Wielandt Theorem (Horn & Johnson, 1985), we have $\left|\frac{1}{n}\lambda_i(K_n) - \frac{1}{n}\lambda_i(K_n^n)\right| \leq \frac{1}{n}\|E\|$. The operator norm of E can be computed as follows.

$$\begin{aligned} \|E\| &= \|K_n - K_n^n\| = \sup_{\mathbf{z} \in \mathcal{S}^{n-1}} |\langle (K_n - K_n^n)\mathbf{z}, \mathbf{z} \rangle| \\ &= \frac{2}{n} \sup_{\mathbf{z} \in \mathcal{S}^{n-1}} \left| \sum_{i=1}^n z_1 z_i [f(\|\mathbf{x}_i - \mathbf{x}_1\|_2^2) - f(\|\mathbf{x}_i - \mathbf{x}_{1'}\|_2^2)] \right| \\ &\leq \frac{2}{n} \sup_{\mathbf{z} \in \mathcal{S}^{n-1}} |z_1| \left| \sum_{i=1}^n z_i [f(\|\mathbf{x}_i - \mathbf{x}_1\|_2^2) - f(\|\mathbf{x}_i - \mathbf{x}_{1'}\|_2^2)] \right| \\ &\leq \frac{2}{n} \sup_{\mathbf{z} \in \mathcal{S}^{n-1}} \left(\sum_{i=1}^n z_i^2 \right)^{\frac{1}{2}} \\ &\quad \cdot \left(\sum_{i=1}^n [f(\|\mathbf{x}_i - \mathbf{x}_1\|_2^2) - f(\|\mathbf{x}_i - \mathbf{x}_{1'}\|_2^2)]^2 \right)^{\frac{1}{2}} \\ &\leq \frac{2|f|_L}{n} \left(\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{x}_1\|_2^2 - \|\mathbf{x}_i - \mathbf{x}_{1'}\|_2^2 \right)^{\frac{1}{2}}, \end{aligned}$$

where we used Hölder's inequality. Now, by the bounded Manifold assumption, i.e. $\|\mathbf{y}_i\| \leq M$, we

have,

$$\begin{aligned} &\|\mathbf{x}_i - \mathbf{x}_1\|_2^2 - \|\mathbf{x}_i - \mathbf{x}_{1'}\|_2^2 \\ &= |\mathbf{x}_1^\top \mathbf{x}_1 - \mathbf{x}_{1'}^\top \mathbf{x}_{1'} - 2\mathbf{x}_1^\top \mathbf{x}_i + 2\mathbf{x}_{1'}^\top \mathbf{x}_i| \\ &= |\mathbf{y}_1^\top \Sigma \mathbf{y}_1 + 2\mathbf{y}_{1'}^\top \Sigma \mathbf{y}_i - \mathbf{y}_1^\top \Sigma \mathbf{y}_{1'} - 2\mathbf{y}_{1'}^\top \Sigma \mathbf{y}_i| \\ &\leq 3M^2 \lambda_{1,p}(\Sigma). \end{aligned}$$

Therefore, the spectral norm of the error matrix can be bounded by,

$$\|E\| \leq \frac{6M^2}{\sqrt{n}} |f|_L \lambda_{1,p}(\Sigma). \quad (3)$$

By the bounded difference inequality and eq. 3 we obtain the desired result. \square

Similar results to Corollary 3.1 can be derived, where the kernel eigenvalues in the exponential bound are replaced by eigenvalues of the sample covariance matrix, in the same way as they appear in the exponential bound of Theorem 3.2. The concentration bounds for the inner product kernel can be derived in a similar way as in Theorem 3.2. The sketch of the derivation is provided in Remark 3.1.

Remark 3.1 (concentration bound for smooth inner product kernels). *For the inner product kernel as defined earlier in this section with $|f|_L$ -Lipschitz function we have $P\left\{\left|\frac{1}{n}\lambda_i(K_n) - \mathbb{E}_S\frac{1}{n}\lambda_i(K_n)\right| > \epsilon\right\} \leq \exp\left\{-\frac{n^2\epsilon^2}{4|f|_L^2 M^4 \lambda_{1,p}^2(\Sigma)}\right\}$. The proof is similar to as of Theorem 3.2. \triangleleft*

Theorem 3.2 and the result in Remark 3.1 connects the concentration of the eigenvalues of the distance and inner product kernel matrices to the difference between the largest and smallest eigenvalues of the covariance matrix. In high dimensional data, the first eigenvalue of the sample covariance matrix is likely to become rather large, implying that the eigenvalues of the distance kernel matrix are not concentrated around the mean. This result suggests that it may be better not to use smooth distance, or inner product kernels in high dimensional data. The results of this section hold also for a wider class of kernel functions that behaves almost like the distance or the inner product kernels. The characterization of the wider class of this functions is summarized in Remark 3.2.

Remark 3.2. *The result of theorem 3.2 also holds for other kernels that satisfy one of the following conditions.*

1. $|k(\mathbf{x}_1, \mathbf{x}_2) - k(\mathbf{x}'_1, \mathbf{x}'_2)| \leq C_{\mathcal{X}} \|\mathbf{x}_1 - \mathbf{x}'_1\|_2^2 - \|\mathbf{x}_2 - \mathbf{x}'_2\|_2^2$,
2. $|k(\mathbf{x}_1, \mathbf{x}_2) - k(\mathbf{x}'_1, \mathbf{x}'_2)| \leq C'_{\mathcal{X}} |\mathbf{x}_1^\top \mathbf{x}'_1 - \mathbf{x}_2^\top \mathbf{x}'_2|$,

where $C_{\mathcal{X}}$ and $C'_{\mathcal{X}}$ are constants depending on the smoothness of the kernel function and the data distribution. \triangleleft

In order to check the bounds in Theorems 3.1 and 3.2, we run two numerical experiments with Gaussian kernels. The details and results are summarized in Example 1.

Example 1. We draw 100 samples from the standard normal distribution, and compute the kernel matrix using a Gaussian kernel function $k(x, y) = \exp(-0.5\|x - y\|_2^2)$. We repeat this procedure a 1000 times, and compute the empirical estimate of the left hand and right hand sides of inequality (1). The results are depicted in Figure 2 (top), where solid, dashed, and dotted lines correspond to left hand side of eq. (1) for $i = 1, 2, 3$, respectively. Lines with different marks correspond to right hand side of inequality (1) for $i = 1, 2, 3$, respectively. As can be seen, the concentration of eigenvalues varies by their order. The bound presented in eq. (1) tracks the concentration changes for each eigenvalue separately.

Similarly we repeat the above example for multivariate Gaussian samples $\mathcal{N}(0, I_2)$ and $\mathcal{N}(0, I_5)$, to check the bound in Theorem 3.2. I_p denotes the identity matrix in \mathbb{R}^p . The empirical estimate of left hand side and right hand side of inequality (2) are drawn in Figure 2 (bottom). As the experiments illustrate, the concentration of eigenvalues change with the dimension of samples, which can be captured by the bound provided in inequality (2).

One of the main obstacles to improve the bounds derived either in this or earlier section is due to the limitation of the matrix perturbation results. For example, the results in Theorem 3.2 and Theorem 2.1 rely on the first order eigenvalue expansion $\lambda_i(\tilde{K}) = \lambda_i(K) - \mathbf{u}_i^\top(K)E\mathbf{u}_i(K) + O(\|K - \tilde{K}\|^2)$, where \tilde{K} is a perturbation of the matrix or compact linear operator K and $E = \tilde{K} - K$. However, these results can be slightly improved by using the second order expansion, as described in (Kato, 1996):

$$\begin{aligned} \lambda_i(\tilde{K}) &= \lambda_i(K) - \mathbf{u}_i^\top(K)E\mathbf{u}_i(K) \\ &+ \sum_{\substack{j=1 \\ j \neq i}}^n \frac{\mathbf{u}_j^\top(K)E\mathbf{u}_i(K)\mathbf{u}_i^\top(K)E\mathbf{u}_j(K)}{\lambda_j - \lambda_i} + O(\|E\|^3). \end{aligned}$$

The above expansion holds if the spectral norm of the error operator/matrix E is smaller than half of the distance between the eigenvalue λ_i and the rest of the spectrum of K . In the notation of Theorem 3.2, $E = K_n - K_n^n$ is symmetric. Therefore, the perturbation of the eigenvalues of the kernel matrix can be bounded by

$$|\lambda_i(K_n) - \lambda_i(K_n^n)| \leq \|E\|_2 + \|E\|_2^2 \sum_{\substack{j=1 \\ j \neq i}}^n \frac{1}{|\lambda_j - \lambda_i|}. \quad (4)$$

Combining the preceding bound with eq. (3) reads

$$P \left\{ \left| \frac{1}{n} \lambda_i(K_n) - \mathbb{E}_S \frac{1}{n} \lambda_i(K_n) \right| > \epsilon \right\} \leq \exp \left\{ -\frac{n^2 \epsilon^2}{\gamma^2} \right\},$$

where

$$\gamma = 6M^2 |f|_L \frac{\lambda_{1,p}(\Sigma)}{\sqrt{n}} + 36M^4 |f|_L^2 \frac{\lambda_{1,p}^2(\Sigma)}{n} \sum_{\substack{j=1 \\ j \neq i}}^n \frac{1}{\lambda_{j,i}^2(K_n)},$$

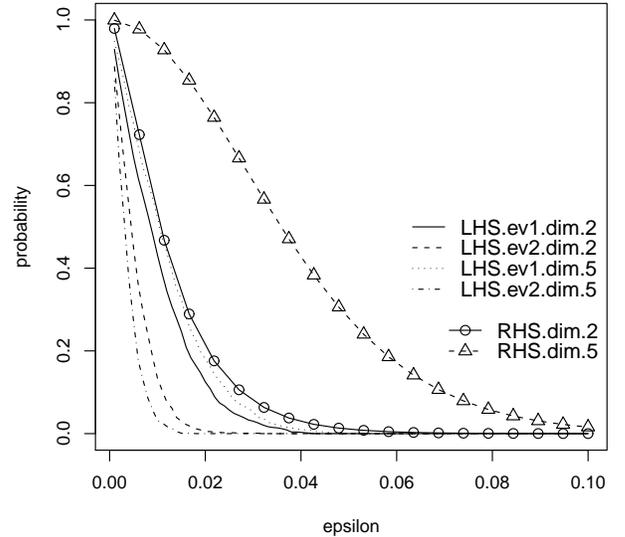
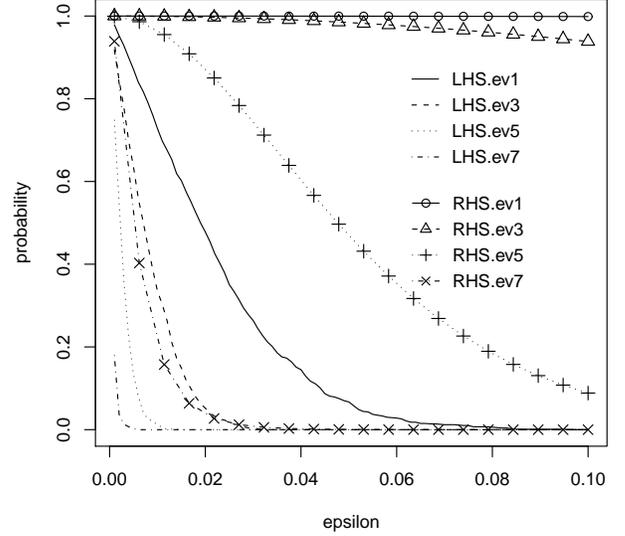


Figure 2: Empirical study on concentrations of eigenvalues. The top figure corresponds to inequality (1) for 1st, 3rd, and 5th eigenvalues. The bottom figure corresponds to inequality (2). See Example 1 for details.

and $\lambda_{j,i}(K_n) = \lambda_j(K_n) - \lambda_i(K_n)$.

Another direction to improve this result is to improve the computation of the spectral norm $\|K_n - K_n^n\|$, so that the dependencies between elements of E are taken into account. This is left for future work.

3.2 EIGENVECTORS OF INNER PRODUCT AND DISTANCE KERNELS

In this section, we provide the concentration bounds for the eigenvectors of the distance and inner product kernel matrices. The results mainly rely on the eigenfunction perturbation expansion summarized in the following theorem, which is due to (Kato, 1996).

Theorem 3.3. [eigenfunction expansion] *Let $\tilde{K} = K + E$, and assume $\|E\|$ is smaller than half of the distance between eigenvalue λ_i and the rest of the eigenvalues, then*

$$\tilde{\mathbf{u}}_i = \mathbf{u}_i + \sum_{\substack{j=1 \\ j \neq i}}^n \frac{\mathbf{u}_j^\top E \mathbf{u}_i}{\lambda_j - \lambda_i} \mathbf{u}_j + O(\|E\|^2), \quad (5)$$

where \mathbf{u}_i is the i -th eigenvector of matrix K , with corresponding eigenvalue λ_i , and $\tilde{\mathbf{u}}$ is the eigenvector of \tilde{K} corresponding to λ_i .

Note that, by eq. (5), when $|\lambda_j - \lambda_i|$ is small, the term $\frac{E}{\lambda_j - \lambda_i}$ has a large norm, which results in an instability of eigenvector \mathbf{u}_i due to the perturbation. This is the case for all eigenvectors that correspond to small eigenvalues. Such phenomenon is empirically studied in (Ng et al., 2001). In the following, we provide pointwise and uniform bounds for the concentration of eigenvectors using eq. (5).

The following results are slightly different from those in (Zwald & Blanchard, 2006), as it contain one extra projection into the given samples. Moreover, the proof presented here holds also for kernels with an infinite dimensional feature space. In the following we denote the evaluation of function $\mathbf{u}(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$ on samples \mathcal{S} by $\mathbf{u}|_{\mathcal{S}} := [\mathbf{u}(\mathbf{x}_1), \dots, \mathbf{u}(\mathbf{x}_n)]$.

Lemma 1. *Let us take the assumptions made in Theorem 3.2. We consider a real valued function $f_{\mathbf{w}}(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto \langle \cdot, \mathbf{w} \rangle$, $\mathbf{w} \in \mathbb{R}^n$. Further, assume $k(\mathbf{x}, \mathbf{y}) = f(\|\mathbf{x} - \mathbf{y}\|_2^2)$. Then, for every $\mathbf{w} \in \mathcal{S}^{n-1}$ we have the following concentration result:*

$$P \{ |f_{\mathbf{w}}(\mathbf{u}_i(K_n)) - \mathbb{E}_{\mathcal{S}} f_{\mathbf{w}}(\mathbf{u}_i(K_n))| > \epsilon \} \leq \exp \left(- \frac{\epsilon^2}{18M^4 |f|_L^2 R_i^2(K_n) \lambda_{1,p}^2(\Sigma)} \right),$$

where $\mathbf{u}_i(K_n)$ is the eigenvector of kernel matrix K_n

corresponding to i -th eigenvalue, $\lambda_{1,p}(\Sigma) := \lambda_1(\Sigma) - \lambda_p(\Sigma)$, and $R_i(K_n) := \sum_{\substack{j=1 \\ j \neq i}}^n \frac{1}{|\lambda_i(K_n) - \lambda_j(K_n)|}$.

Proof. Let us define the perturbed matrix K_n^n by $[K_n^n]_{i,j} = \frac{1}{n} k(\mathbf{x}_i, \mathbf{x}_j)$, $\forall i, j = 2, \dots, n$ and $[K_n^n]_{1,j} = K_{j,1}^n = k(\mathbf{x}'_1, \mathbf{x}_j)$. Therefore, $K_n = K_n^n + E$. By definition of $f_{\mathbf{w}}$ we have,

$$|f_{\mathbf{w}}(\mathbf{u}_i(K_n)) - f_{\mathbf{w}}(\mathbf{u}_i(K_n^n))| = |f_{\mathbf{w}}(\mathbf{u}_i(K_n) - \mathbf{u}_i(K_n^n))| \leq \|\mathbf{w}\| \|\mathbf{u}_i(K_n) - \mathbf{u}_i(K_n^n)\|,$$

where $\mathbf{u}_i(K_n)$ and $\mathbf{u}_i(K_n^n)$ are eigenvectors of matrices K_n and K_n^n . Now, we need to compute the quantity $\|\mathbf{u}_i(K_n) - \mathbf{u}_i(K_n^n)\|$. Using the perturbation result of Theorem 3.3, we have $\|\mathbf{u}_i(K_n) - \mathbf{u}_i(K_n^n)\| \leq \|E\| R_i(K_n)$. Then, for all $\mathbf{w} \in \mathcal{S}^{n-1}$, applying the bounded difference inequality result in the following concentration bound.

$$P \{ |\langle \mathbf{w}, \mathbf{u}_i(K_n) - \mathbb{E}_{\mathcal{S}} \mathbf{u}_i(K_n) \rangle| > \epsilon \} \leq \exp \left(- \frac{2\epsilon^2}{n \|E\|^2 R_i^2(K_n) \sup_{\mathbf{w}} \|\mathbf{w}\|_2} \right). \quad (6)$$

By plugging eq. (3) into eq. (6) we obtain the claimed result. \square

The concentration bound in Lemma 1 can be extended to real-valued convex Lipschitz functions. An immediate application of Lemma 1 is to derive a uniform concentration bound for the kernel matrix eigenvectors. This result is provided in Corollary 3.2.

Corollary 3.2 (uniform concentration of kernel eigenvectors). *With the same assumptions made in Theorem 1, we have*

$$P \{ \|\mathbf{u}_i(K_n) - \mathbb{E}_{\mathcal{S}} \mathbf{u}_i(K_n)|_{\mathcal{S}}\| > \epsilon \} \leq 2 \exp(2n - c\epsilon^2),$$

where $c^{-1} = 18M^4 |f|_L^2 R_i^2(K_n) \lambda_{1,p}^2(\Sigma)$.

In the above bound, the positive first term inside the exponential function does not bring any harm, since the first eigenvalue of the sample covariance matrix is of order $O(\sqrt{n})$, which cancels out the effect of the $2n$ term inside the exponential function. The proof relies on Lemma 1 and the ε -Net argument (Ledoux & Talagrand, 1991).

Proof of Corollary 3.2. Let us denote the ε -Net of the compact unit sphere \mathcal{S}^{n-1} by $\mathcal{N}_{\varepsilon}$, for some $\varepsilon \in (0, 1)$. Let $\mathbf{x} \in \mathcal{S}^{n-1}$ such that $\|\mathbf{u}\| = \langle \mathbf{x}, \mathbf{u} \rangle$, where $\mathbf{u} := \mathbf{u}_i(K_n)$. Now, we can choose $\mathbf{y} \in \mathcal{N}_{\varepsilon}$ so that $\|\mathbf{x} - \mathbf{y}\| \leq \varepsilon$ holds. Then, we have $|\langle \mathbf{x}, \mathbf{u} \rangle - \langle \mathbf{y}, \mathbf{u} \rangle| \leq \|\mathbf{u}\| \|\mathbf{x} - \mathbf{y}\| \leq \varepsilon \|\mathbf{u}\|$. By the triangle inequality, we have $|\langle \mathbf{y}, \mathbf{u} \rangle| \geq |\langle \mathbf{x}, \mathbf{u} \rangle| - |\langle \mathbf{x}, \mathbf{u} \rangle - \langle \mathbf{y}, \mathbf{u} \rangle| \geq \|\mathbf{u}\| - \varepsilon \|\mathbf{u}\|$. Putting all together, we obtain

$$\max_{\mathbf{y} \in \mathcal{N}_{\varepsilon}} |\langle \mathbf{y}, \mathbf{u} \rangle| \leq \|\mathbf{u}\| \leq \frac{1}{1 - \varepsilon} \max_{\mathbf{y} \in \mathcal{N}_{\varepsilon}} |\langle \mathbf{y}, \mathbf{u} \rangle|.$$

By using the upper bound in the preceding inequality, we have

$$\begin{aligned}
& P\{\|\mathbf{u} - \mathbb{E}_{\mathcal{S}}\mathbf{u}|_{\mathcal{S}}\| > \epsilon\} \\
&= P\left\{\sup_{\mathbf{w} \in \mathcal{S}^{n-1}} |\langle \mathbf{w}, \mathbf{u} - \mathbb{E}_{\mathcal{S}}\mathbf{u}|_{\mathcal{S}} \rangle| > \epsilon\right\} \\
&\leq P\left\{\frac{1}{1-\epsilon} \max_{\mathbf{w} \in \mathcal{N}_{\epsilon}} |\langle \mathbf{w}, \mathbf{u} - \mathbb{E}_{\mathcal{S}}\mathbf{u}|_{\mathcal{S}} \rangle| > \epsilon\right\} \\
&\leq |\mathcal{N}_{\epsilon}| P\{|\langle \mathbf{w}, \mathbf{u} - \mathbb{E}_{\mathcal{S}}\mathbf{u}|_{\mathcal{S}} \rangle| > \epsilon\} \\
&\leq 2|\mathcal{N}_{\epsilon}| \exp(-c\epsilon^2),
\end{aligned}$$

where in the second line we have used the union bound. By Lemma 9.5. in (Ledoux & Talagrand, 1991), for $\epsilon = \frac{1}{2}$ we have $|\mathcal{N}_{\epsilon}| \leq 6^n$. Putting all together we get the claimed result. \square

The result on the eigenvectors of distance kernel can be applied to a broader class of kernels that satisfies Lipschitz smoothness, see Remark 3.2. In particular, for the inner product kernels, we have $\|E\| \leq 2M^2 \frac{|f|_{L^{\lambda_{1,p}}(\Sigma)}}{\sqrt{n}}$. Plugging this into inequality (6), we get the desired bounds.

The concentration bounds presented in Theorem 1 and Corollary 3.2 can also be slightly improved by using a higher order expansion of the operator perturbation. For more details on higher order expansions of eigenvectors see (Kato, 1996).

4 CASE STUDY: KERNEL TARGET-ALIGNMENT

In this section, we provide an example of using the ingredients of the proof in Section 3, for deriving a concentration inequality for the sample kernel alignment. This is a good example to show how the combination of a proper matrix perturbation result and the concentration inequality provides an informative, simple and easy to compute concentration bound.

Kernel target-alignment is proposed in (Cristianini et al., 2002), for measuring the agreement between a kernel matrix and the given learning task, i.e. classification. By changing the eigenvalues of the kernel matrix, one can improve the alignment to the labels, i.e. target values, which results in improvement of the learning performance. This approach has been proposed for kernel learning (for more detail see (Cristianini et al., 2002)). (Jia & Liao, 2009) proposed to use the bound in Theorem 2.2 to derive a concentration inequality for sample kernel alignment. Suppose a sample set $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \{-1, +1\}$ is given. Let $K_n \in \mathbb{R}^{n \times n}$ be a Mercer kernel matrix defined by $[K_n]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j), \forall 1 \leq i, j \leq n$. The kernel

target-alignment or kernel alignment is defined by

$$A(y) = \frac{\langle \mathbf{y} \otimes \mathbf{y}, k \rangle}{\|\mathbf{y} \otimes \mathbf{y}\| \|k\|} = \frac{\langle \mathbf{y} \otimes \mathbf{y}, k \rangle}{\|k\|},$$

where $\langle f, g \rangle = \int_{\mathcal{X}^2} f(x, z)g(x, z)dP(x)dP(z)$. Similarly, the sample kernel target alignment of K_n is defined by

$$A(K_n) = \frac{\langle K_n, Y \otimes Y \rangle_F}{\sqrt{\langle Y \otimes Y, Y \otimes Y \rangle_F \langle K, K_n \rangle}} = \frac{Y^\top K_n Y}{n \|K_n\|_F},$$

where $Y = (y_1, \dots, y_n)$, $\langle \cdot, \cdot \rangle_F$ is the Frobenius inner product and $\|\cdot\|_F$ is the Frobenius norm. In this section, we drop the subscript F . The following bound is suggested in (Jia & Liao, 2009),

$$P\{|A(K_n) - A(y)| > \epsilon\} \leq 2 \exp\left(-\frac{2\epsilon^2(n-1)^2}{nC^2(\theta)}\right),$$

where $C(\theta) = |A(K_n)|\theta^{-1}(m - (m-1)\theta + \frac{2n-1}{\|K_n\|})$ and θ is defined by $\theta = 1 - \max_{s=1, \dots, n} \min_{i=1, \dots, n-1} \frac{\lambda_i(K_n^s)}{\lambda_i(K_n)}$, where K_n^s is the submatrix of K_n , where the s -th row and column are replaced by $\mathbf{0}$ vectors. Here, we provide another concentration bound, which depends only on eigenvalues of the sample kernel matrix, and can be approximated by the ratio between the first and second largest eigenvalues of the kernel matrix.

Theorem 4.1. *Let us take (Assumption 1), and $A(y)$ and $A(K)$ are defined as above. Then, we have*

$$\begin{aligned}
& P\{|A(K_n) - A(y)| > \epsilon\} \\
&\leq 2 \exp\left(-\frac{2\epsilon^2}{A(K_n) \left|\frac{1}{n-1} - \frac{\|K_n\|}{L}\right| + \left(2 + \frac{1}{n-1}\right) \frac{1}{L}}\right).
\end{aligned}$$

Proof.(Sketch) The proof follows from the Cauchy interlacing lemma and the bounded difference inequality. Let us define K_n^n as K_n^s as above, where $s = n$. Then, we have

$$\begin{aligned}
|A(K_n) - A(K_n^n)| &\leq \left| \frac{\langle K_n, YY^\top \rangle}{n \|K_n\|} - \frac{\langle K_n^n, Y^n Y^{n\top} \rangle}{(n-1) \|K_n^n\|} \right| \\
&\leq \left| \frac{\langle K_n, YY^\top \rangle}{n \|K_n\|} - \frac{\langle K_n, YY^\top \rangle}{(n-1) \|K_n^n\|} \right| \\
&\quad + \left| \frac{\langle K_n, YY^\top \rangle}{(n-1) \|K_n^n\|} - \frac{\langle K_n^n, Y^n Y^{n\top} \rangle}{(n-1) \|K_n^n\|} \right| \\
&\leq |\langle K_n, YY^\top \rangle_F| \left| \frac{1}{(n-1) \|K_n^n\|} - \frac{1}{\|K_n\|} \right| \\
&\quad + \frac{2n-1}{(n-1) \|K_n^n\|}.
\end{aligned}$$

The result follows by plugging in the following inequality into the above.

$$L := \sqrt{\sum_{i=2}^{n-1} \lambda_i^2(K_n)} \leq \|K_n^n\|_F \leq \sqrt{\sum_{i=1}^{n-1} \lambda_i^2(K_n)}.$$

Therefore, we have

$$\begin{aligned} & |A(K_n) - A(K_n^n)| \\ &= \frac{\langle K_n, YY^\top \rangle}{\|K_n\|} \left| \frac{(n-1)\|K_n\| - \|K_n^n\|}{(n-1)\|K_n^n\|} \right| + \frac{2n-1}{n-1} \frac{1}{\|K_n^n\|} \\ &\leq A(K_n) \left| \frac{1}{n-1} - \frac{\|K_n\|}{L} \right| + \left(2 + \frac{1}{n-1} \right) \frac{1}{L}. \end{aligned}$$

□

In the preceding results, the term $\frac{\|K_n\|}{L}$ can be approximated by $\frac{\lambda_1(K_n)}{\lambda_2(K_n)}$. Therefore, for sufficiently large samples size, the gap between first and second eigenvalues control the concentration of the sample kernel alignment.

5 CONCLUDING REMARKS

Kernel methods are very successful approaches in solving different machine learning problems. This success is mainly rooted in using feature maps and kernel functions, which transform the given samples into a possibly higher dimensional space.

Kernel matrices/operators can be characterized by the properties of their spectral decomposition. Also, in the computation of the excess risk, we eventually need the eigenvalue information. Concentration inequalities for the spectrum of kernel matrices measure how close the sample eigenvalues/eigenvectors are to the population values. The previous concentration results are either suboptimal, or computing the bound is impractical. This paper improves upon the optimality and computational feasibility of the previous results, by applying Cauchy's interlacing lemma. Moreover, for inner product and Euclidean distance kernels, e.g. radial basis function (Steinwart & Christmann, 2008), we derive new types of bounds, which connect the concentration of sample kernel eigenvalues and eigenvectors to the eigenvalues of the sample covariance matrix. This result may explain the poor performance of some nonlinear kernels in very high dimensions, as the largest eigenvalue of the sample covariance matrix become very large. As an interesting case study, we establish a computationally less demanding, simple and informative bound for the concentration of the sample kernel target-alignment.

The results can be improved, for example by using more careful calculations of the spectral norms in bounded difference computation, other operator perturbation results, or local concentration inequalities.

Acknowledgements

Authors acknowledge the reviewers' comments, which improved the paper considerably. NR and RV were

funded by the Academy of Finland, through Centres of Excellence Program 2006-2011.

References

- [1] Bengio, Y., Vincent, P., Paiement, J., Delalleau, O., Ouimet, M., Le Roux, N., Spectral clustering and kernel PCA are learning eigenfunctions. Technical Report TR 1239, University of Montreal, 2003.
- [2] Jia, L. and Liao, S., Accurate probabilistic error bound for eigenvalues of kernel matrix, LNCS, Vol. 5828/2009, 162-175, 2009.
- [3] Mendelson, S., and Pajor, A., Ellipsoid approximation using random vectors, COLT 2005, Springer-Verlag, LNAI 3359, 429-443, 2005.
- [4] Koltchinskii, V., Asymptotics of spectral projections of some random matrices approximating integral operators. Progress in Probability 43, 191-227, 1998.
- [5] Koltchinskii, V., Giné, E., Random matrix approximation of spectra of integral operators, Bernoulli 6(1), 113-167, 2000.
- [6] Shawe-Taylor, J., Williams, C.K., Cristianini, N., Kandola, J., On the eigenspectrum of the gram matrix and the generalization error of kernel-PCA. IEEE Transactions on Information Theory 51(7), 2510-2522, 2005.
- [7] Zwald, L., Blanchard, G., On the convergence of eigenspaces in kernel principal component analysis, NIPS,1649-1656, MIT Press, Cambridge, 2006.
- [8] Cristianini, N., Kandola, J., Elisseeff, A., Shawe-Taylor, J., On kernel-target alignment, Journal of Machine Learning Research 1, 131,2002.
- [9] Steinwart, I., and Christmann, A., *Support Vector Machines*, Springer, 2008.
- [10] Horn, R. , Johnson, C., Matrix Analysis, Cambridge University Press, New York, 1985.
- [11] Andrew Y. Ng , Alice X. Zheng , Michael I. Jordan, Link analysis, eigenvectors and stability, Proceedings of the 17th international joint conference on Artificial Intelligence, 903-910, Seattle, USA, 2001.
- [12] Kato, T., *Perturbation theory for linear operators*, Springer-Verlag, New York, 1996.
- [13] Ledoux, M., and Talagrand, M., *Probability in Banach Spaces*, Isoperimetry and processes, Springer-Verlag, Berlin, 1991.