

# Davis-Kahan Theorem

Yutong Wang  
London School of Economics

Mar 2024

## Contents

<b>1</b>	<b>Principal Angles</b>	<b>2</b>
<b>2</b>	<b>Singular Value Decomposition</b>	<b>3</b>
<b>3</b>	<b>Norm</b>	<b>3</b>
3.1	Norm of vectors . . . . .	4
3.2	Norm of matrices . . . . .	4
<b>4</b>	<b>Projection</b>	<b>5</b>
<b>5</b>	<b>Inequalities</b>	<b>5</b>
<b>6</b>	<b>Davis-Kahan Theorem</b>	<b>6</b>
<b>7</b>	<b>A Useful Variant</b>	<b>8</b>

Davis-Kahan Theorem is about how to give a bound of the eigenspace or eigenvectors after the original matrix is perturbed. To understand it, basic knowledge of linear algebra is required.

Here we focus on the Symmetric matrix. Assume  $A \in \mathbb{R}^{n \times n}$  is symmetric, where  $A+H \in \mathbb{R}^{n \times n}$  is still symmetric (which implies the perturbation  $H$  is still symmetric). Thus we have the eigendecomposition:

$$\begin{aligned} A &= \sum_{i=1}^n \lambda_i u_i u_i^*, \\ A+H &= \sum_{i=1}^n \mu_i v_i v_i^*. \end{aligned}$$

We would like to show that if  $H$  is ‘small’ (in some sense), then  $u_i$  is close to  $v_i$ . One of the key concerns is that the order of the spectrum may change after the perturbation. See a typical example:

**Example 1.** Let  $A = \begin{pmatrix} 1 & \varepsilon \\ \varepsilon & 1 \end{pmatrix}$ , whose eigenvalues are the roots of equation  $(1-\lambda)^2 - \varepsilon^2 = 0$ , i.e.,  $\lambda_1 = 1 + \varepsilon$ ,  $\lambda_2 = 1 - \varepsilon$ .

## 1 Principal Angles

The Principal angles are defined step by step. First, we could define the principal angle between two vectors; then the principal angle between two lines (lines are 1-D subspaces) can be defined; Next, the definition can be generalised to principal angles between a line and a 2-D subspace; and eventually principal angles between any two sub-spaces (no further assumptions on the dimension of two sub-spaces).

From the process of defining, we could see that the principal angle can be defined in any inner-product space. Now let’s consider Euclidean spaces only.

1. **Between vectors.** Assume  $x, y \in \mathbb{R}^n$ , then the principal angle between them is defined via their as:

$$\theta_{x,y} := \arccos \frac{x^\top y}{\|x\| \|y\|}.$$

It is worth mentioning that the RHS can take any value between  $[-1, 1]$ . To cover the range, we allow  $\theta_{x,y}$  taking value in  $[0, \pi]$ . It is easy to construct intuition by setting  $n = 2$ , where all the stuff can be drawn in a 2-D paper.

2. **Between lines.** Assume still  $x, y \in \mathbb{R}^n$ , We are going to define the principal angle for two subspaces:  $\mathbb{R}x := \{ax \mid a \in \mathbb{R}\}$  and  $\mathbb{R}y := \{ay \mid a \in \mathbb{R}\}$ .

$$\theta_{\mathbb{R}x, \mathbb{R}y} := \arccos \frac{|x^\top y|}{\|x\| \|y\|} = \inf_{a \in \mathbb{R}, b \in \mathbb{R}} \arccos \frac{|(ax)^\top (by)|}{\|ax\| \|by\|}.$$

Here the infimum does not actually solve any optimisation, since the objective function does not depend on either  $a$  or  $b$ . However, this expression is convenient for generalisation.

In addition, one can see that  $\cos(\theta_{\mathbb{R}x, \mathbb{R}y})$  is always non-negative, thus, allowing  $\theta_{\mathbb{R}x, \mathbb{R}y} \in [0, \frac{\pi}{2}]$  is enough to cover all values.

3. **Between a line and a plane (2-D subspace).** Now assume we have a line spanned by  $x \in \mathbb{R}^n$ , i.e.,  $\mathbb{R}x = \{ax \mid a \in \mathbb{R}\}$ , and a 2-D subspace spanned by vector  $y, z \in \mathbb{R}^n$ , i.e.  $\mathbb{R}\{y, z\} := \{by + cz \mid b, c \in \mathbb{R}\}$ . The principal angle between the line and the plane is:

$$\theta_{\mathbb{R}x, \mathbb{R}\{y, z\}} := \inf_{a, b, c \in \mathbb{R}} \arccos \frac{|(ax)^\top (by + cz)|}{\|ax\| \|by + cz\|}.$$

4. **Between any two subspace of  $\mathbb{R}^n$ .** Assume the two subspaces  $E, F$  are spanned by  $\{e_1, \dots, e_l\}$  and  $\{f_1, \dots, f_m\}$ , respectively. WLOG, we could assume  $E = (e_1, \dots, e_l) \in \mathbb{R}^{n \times l}$  and  $F = (f_1, \dots, f_m) \in \mathbb{R}^{n \times m}$  are matrices consist of two orthonormal basis.<sup>12</sup> By the definition of the basis, we know that any vector in space  $E$  can be expressed as  $E\alpha$  for  $\alpha \in \mathbb{R}^l$ , any vector in space  $F$  can be expressed as  $F\beta$  for  $\beta \in \mathbb{R}^m$ , the principal angle between the two subspaces  $E, F$  is defined as:

$$\theta_{E,F} = \inf_{\|\alpha\|=1, \|\beta\|=1} \arccos \left( (E\alpha)^\top (F\beta) \right) = \inf_{\|\alpha\|=1, \|\beta\|=1} \arccos \left( \alpha^\top E^\top F \beta \right)$$

Now the definition of principal angles has been clear, we could apply several tools to analyse it.

## 2 Singular Value Decomposition

One characterisation of singular value is by the min-max Theorem.

**Theorem 1.** (*Min-max Theorem*) For a matrix  $A \in \mathbb{R}^{m \times n}$ , it has  $\min\{m, n\}$  number of singular values, and they can be characterised as: for  $i = 1, \dots, \min\{m, n\}$ ,

$$\sigma_i(A) = \min_{\dim(U)=n-i+1} \max_{x \in U, \|x\|_2=1} \|Ax\|_2$$

The importance of singular value partly lies in its unitary invariance.

**Theorem 2.** (*Unitary invariance*) For unitary matrix  $U$  and  $V$  (which makes all the matrix multiples compatible), the singular value of  $A$  and  $UAV$  are the same.

Combine the above two results, one can see that, for a symmetric (and of course, squared) matrix  $A$ , its singular value can be characterised as:

$$\sigma_i(A) = \min_{\dim(U)=n-i+1} \max_{x \in U, \|x\|_2=1, y \in U, \|y\|_2=1} \|y^\top Ax\|_2.$$

With this in mind, let's take a look at the definition of the principal angle between two subspaces. Since  $\arccos(\cdot)$  is a decreasing function, to find the infimum, we want what is inside  $\arccos$  as large as possible, thus:

$$\theta_{E,F} = \arccos \sigma_{\max}(E^\top F).$$

In this way, we relate the geometry (the angle) to the algebra (the singular value). To study all angles rather than the 'principal' angle, we utilise all singular values of  $E^\top F$ . Note that  $E^\top F \in \mathbb{R}^{l \times m}$ , the singular values can be expressed as  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k$ , where  $k = \min\{l, m\}$ . Define the angle  $\Theta_{E,F}$  between space  $E$  and  $F$  as:

$$\cos(\Theta_{E,F}) = \text{diag}(\sigma_1, \dots, \sigma_k).$$

## 3 Norm

Normed space is a linear space equipped with a structure called a norm. Intuitively, norm is a way of defining 'length'.

**Definition 1.** Given a vector space  $V$ , a real-valued function  $p : V \rightarrow \mathbb{R}$  is called a norm if:

<sup>1</sup>Check. It is a useful practice to get familiar with linear algebra.

<sup>2</sup>There is a bit of abuse of notations:  $E$  denotes both the basis and the subspace.

1. (Sub-additivity)  $p(x + y) \leq p(x) + p(y)$  for all  $x, y \in V$ .
2. (Absolute homogeneity)  $p(sx) = |s|p(x)$  for any scalar  $s$ .
3. (Positive definiteness) If  $p(x) = 0$ , then  $x = \mathbf{0}$ .

It is worth mentioning that, in a normed space  $(V, \|\cdot\|)$ , there is a naturally induced metric defined as: for any  $x, y \in V$ ,

$$\rho(x, y) := \|x - y\|$$

### 3.1 Norm of vectors

For  $x \in \mathbb{R}^n$ ,

1.  $\|x\|_2 = \left(\sum_{i=1}^n x_i^2\right)^{\frac{1}{2}}$  is a norm, called  $l_2$ -norm.
2.  $\|x\|_1 = \sum_{i=1}^n |x_i|$  is a norm, called  $l_1$ -norm.
3. Generally,  $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{\frac{1}{p}}$  is a norm (for  $p \geq 1$ ), called  $l_p$ -norm.
4. As another extreme case, set  $p \rightarrow \infty$ , one has:

$$\|x\|_\infty = \lim_{p \rightarrow \infty} \left(\sum_{i=1}^n |x_i|^p\right)^{\frac{1}{p}} = \max_{i=1, \dots, n} |x_i|,$$

called the  $l_\infty$ -norm

### 3.2 Norm of matrices

One way to look at a  $m \times n$  matrix is to treat it as a linear operator from  $\mathbb{R}^n$  to  $\mathbb{R}^m$  under a particular basis. In this sense, the norm of matrices can be defined in the following way:

Suppose we have equipped  $\mathbb{R}^n$  with norm  $\|\cdot\|_U$ , and equipped  $\mathbb{R}^m$  with norm  $\|\cdot\|_V$ . Then, for any  $A \in \mathbb{R}^{m \times n}$ , there is a natural induced norm<sup>3</sup>:

$$\|A\| := \sup_{x \in \mathbb{R}^n} \frac{\|Ax\|_V}{\|x\|_U} = \sup_{x \in \mathbb{R}^n, \|x\|_U=1} \|Ax\|_V.$$

A lot of norms of this kind can be defined:

1. The  $2 \rightarrow 2$  norm. Take the  $l_2$  norm on  $\mathbb{R}^m$ , and also  $l_2$  norm on  $\mathbb{R}^n$ . We have:

$$\|A\| := \sup_{x \in \mathbb{R}^n} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{x \in \mathbb{R}^n, \|x\|_2=1} \|Ax\|_2.$$

Combined with the min-max Theorem, one can see that the  $2 \rightarrow 2$  norm is equal to the maximum among absolute values of all singular values.

Another important and widely-used norm is the Frobenius norm, which is quite similar to the  $l_2$ -norm for vector if we stretch the matrix to be a vector in  $\mathbb{R}^{mn \times 1}$ :

---

<sup>3</sup>Verifying why it satisfies all three conditions is a practice to get familiar with the definition.

**Definition 2.** For  $A \in \mathbb{R}^{m \times n}$ , its Frobenius norm is defined by:

$$\|A\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{\frac{1}{2}}.$$

Several important results regarding the Frobenius norm are required.

**Theorem 3.** For  $A \in \mathbb{R}^{m \times n}$ ,

$$\|A\|_F^2 = \text{tr}(AA^\top).$$

## 4 Projection

For any  $k$ -dimensional subspace of  $\mathbb{R}^n$  ( $k \leq n$ ), we could assume it is spanned by an orthonormal basis. This is due to an algorithm to find orthonormal basis for any subspace.

**Theorem 4.** (Gram-Schmidt algorithm) Google it.

**Definition 3.** A function mapping from a vector space  $V$  to itself,  $P : V \rightarrow V$ , is called a projection if:

$$P \circ P = P.$$

**Theorem 5.** If a subspace of  $\mathbb{R}^n$  is spanned by an orthonormal basis  $\{u_1, \dots, u_k\}$ , then the matrix expression of the projection operator is given by:

$$P_u = \sum_{i=1}^k u_i u_i^\top.$$

*Proof.* The sketch:

1. Verifying  $P_u \circ P_u = P_u$  to show it is a projection.
2. Show for any  $u_i$ ,  $i = 1, \dots, k$ ,  $P_u u_i = u_i$ .
3. Show for any vector  $v$  linearly independent from  $\{u_1, \dots, u_k\}$ ,  $P_u v = 0$ .

□

## 5 Inequalities

**Proposition 1.** For diagonal matrix  $D = \text{diag}(d_1, \dots, d_n)$  with diagonal elements lies between  $[d_{\min}, d_{\max}]$ , we have:

$$d_{\min} \|X\|_F \leq \|DX\|_F \leq d_{\max} \|X\|_F$$

*Proof.* Sketch: Treat  $\|DX\|_F^2$  as a function of  $(d_1, \dots, d_n)$ , Show that it is increasing in all  $d_i$ .

$$\|DX\|_F^2 = \text{tr}(DX(DX)^\top)$$

□

## 6 Davis-Kahan Theorem

Now we are ready to talk about the Davis-Kahan Theorem. Recall the setting and the problem that we concern: assume we have a symmetric matrix  $A$ , and a symmetric perturbation  $H$ , with eigendecomposition:

$$\begin{aligned} A &= \sum_{i=1}^n \lambda_i u_i u_i^* = U \Lambda U^\top + U_\perp \Lambda_\perp U_\perp^\top, \\ A + H &= \sum_{i=1}^n \mu_i v_i v_i^* = V M V^\top + V_\perp M_\perp V_\perp^\top. \end{aligned}$$

With the belief that if  $H$  is small, then  $u_i$  and  $v_i$  should be very close. We are interested in quantitatively measuring how close  $u_i$  and  $v_i$  are, and relate it to some quantity of the perturbation  $H$ . Now it is clear, since  $u_i, v_i \in \mathbb{R}^n$ , and  $\|u_i\|_2 = 1, \|v_i\|_2 = 1$ , they are in a unit sphere:  $u_i, v_i \in \mathcal{S}^{n-1}$ . It is reasonable to measure how close they are by the principal angle between the two unit vectors. Furthermore, it is equivalent to measuring angles between the spaces generated by the first  $k$  eigenvectors  $\{u_1, \dots, u_k\}$  and  $\{v_1, \dots, v_k\}$ .

Denote  $U = (u_1, \dots, u_k)$  and  $V = (v_1, \dots, v_k)$ , then  $P := UU^\top$  and  $Q := VV^\top$  are the projection onto space  $U$  and  $V$ .

$$\begin{aligned} \|P - Q\|_F^2 &= \text{tr}((P - Q)(P - Q)^\top) \\ &= \text{tr}(PP^\top - QP^\top - PQ^\top - QQ^\top) \\ &= 2k - 2\text{tr}(PQ^\top) \\ &= 2k - 2\text{tr}(UU^\top(VV^\top)^\top) \\ &= 2k - 2\text{tr}(U(U^\top V V^\top)) \\ &= 2k - 2\text{tr}((U^\top V V^\top)U) \\ &= 2k - 2\text{tr}((U^\top V)(U^\top V)^\top) \\ &= 2k - 2\|U^\top V\|_F^2, \end{aligned}$$

Since the Frobenius norm is unitary invariant, denote the eigendecomposition of  $U^\top V$  as  $U^\top V = W_1 \Sigma W_2^\top$ , we have:

$$\begin{aligned} \|U^\top V\|_F^2 &= \|W_1 \Sigma W_2^\top\|_F^2 \\ &= \|\Sigma\|_F^2 \\ &= \|\cos(\Theta_{U,V})\|_F^2 \\ &= \sum_{i=1}^k \cos^2(\theta_{U,V}). \end{aligned}$$

Combine the results above, we have:

$$\begin{aligned} \|P - Q\|_F^2 &= 2k - 2 \sum_{i=1}^k \cos^2(\theta_{U,V}) \\ &= 2 \sum_{i=1}^k \sin^2(\theta_{U,V}) \\ &= 2\|\sin(\Theta_{U,V})\|_F^2. \end{aligned}$$

Several properties of the function  $\text{tr}(\cdot)$  are used in the derivation, which is summarised in the following proposition.

**Proposition 2.** *For a matrix  $A \in \mathbb{R}^{n \times m}$ , the function  $\text{tr}(\cdot) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  is defined as:*

$$\text{tr}(A) = \sum_{i=1}^{\min\{m,n\}} A_{ii}.$$

The function has the following properties:

1. It is a linear function, i.e., for any  $A, B$ ,  $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$ .
2. For compatible matrices  $A, B$ ,  $\text{tr}(AB) = \text{tr}(BA)$ .

We have shown how to measure the difference between  $u_i$  and  $v_i$ , now we are ready to upper bound the  $\|U^\top V\|_F$ .

Recall that we have the eigendecomposition of  $A$  and  $A + H$  in matrix form,

$$\begin{aligned} A &= \sum_{i=1}^n \lambda_i u_i u_i^* = U \Lambda U^\top + U_\perp \Lambda_\perp U_\perp^\top, \\ A + H &= \sum_{i=1}^n \mu_i v_i v_i^* = V M V^\top + V_\perp M_\perp V_\perp^\top. \end{aligned}$$

Multiply  $U$  on the right of  $H = (A + H) - A$ , we get:

$$\begin{aligned} HU &= (A + H)U - AU \\ &= (V M V^\top + V_\perp M_\perp V_\perp^\top)U - (U \Lambda U^\top + U_\perp \Lambda_\perp U_\perp^\top)U \\ &= (V M V^\top)U + (V_\perp M_\perp V_\perp^\top)U - U \Lambda. \end{aligned}$$

Then, multiply  $V_\perp^\top$  on the left, we get:

$$\begin{aligned} V_\perp^\top HU &= V_\perp^\top \left( (V M V^\top)U + (V_\perp M_\perp V_\perp^\top)U - U \Lambda \right) \\ &= M_\perp V_\perp^\top U - V_\perp^\top U \Lambda. \end{aligned} \tag{6.1}$$

Using the 6.1 together with some inequalities, we could prove the results.

$$\begin{aligned} \|V_\perp^\top HU\|_F &= \|M_\perp V_\perp^\top U - V_\perp^\top U \Lambda\|_F \\ &\geq | \|M_\perp V_\perp^\top U\|_F - \|V_\perp^\top U \Lambda\|_F | \\ &\geq \max\{ \|M_\perp V_\perp^\top U\|_F - \|V_\perp^\top U \Lambda\|_F, \|V_\perp^\top U \Lambda\|_F - \|M_\perp V_\perp^\top U\|_F \} \end{aligned}$$

The last step is because we are not sure which is larger,  $\|V_\perp^\top U \Lambda\|_F$  or  $\|M_\perp V_\perp^\top U\|_F$ . Take the first situation as an example, using the inequality,

$$\|M_\perp V_\perp^\top U\|_F - \|V_\perp^\top U \Lambda\|_F \geq \min\{M_\perp\} \|V_\perp^\top U\|_F - \|V_\perp^\top U\|_F \max\{\Lambda\}$$

From intuition level, we have seen why we need the eigengap for the two matrices  $A$  and  $A + H$ , now we can see from the mathematical expression.

**Assumption 1.** (Eigengap) *There exists a constant  $\delta > 0$ , such that  $|\lambda_i - \mu_j| \geq \delta$  for all  $i = 1, \dots, k$ ,  $j = k + 1, \dots, n$ .*

With the assumption above, we have:

$$\|V_\perp^\top HU\|_F \geq \delta \|V_\perp^\top U\|_F.$$

For the LHS, note that:

$$\begin{aligned} \|V_\perp^\top HU\|_F^2 &= \text{tr}\left((V_\perp^\top HU)(V_\perp^\top HU)^\top\right) \\ &= \text{tr}\left((V_\perp^\top HU)U^\top H^\top V_\perp\right) \\ &= \text{tr}\left(V_\perp^\top H P P H^\top V_\perp\right) \\ &= \text{tr}\left(P H^\top V_\perp V_\perp^\top H P\right) \\ &= \text{tr}\left(P H^\top Q Q H P\right) \\ &\leq \text{tr}(H^\top H) \\ &= \|H\|_F^2. \end{aligned}$$

For the RHS,

$$\begin{aligned} \|V_\perp^\top U\|_F^2 &= \text{tr}\left(V_\perp^\top U (V_\perp^\top U)^\top\right) \\ &= \text{tr}\left(V_\perp^\top U U^\top V_\perp\right) \\ &= \text{tr}\left(U U^\top V_\perp V_\perp^\top\right) \\ &= \text{tr}\left(P(I_n - Q)\right) \\ &= \text{tr}(P) - \text{tr}(PQ) \\ &= k - \text{tr}(PQ) \\ &= k - \text{tr}(U U^\top V V^\top) \\ &= k - \text{tr}(U^\top V V^\top U) \\ &= k - \|\cos(\Theta_{U,V})\|_F^2 \\ &= \|\sin(\Theta_{U,V})\|_F^2. \end{aligned}$$

Combine all the inequalities, we get:

$$\|\sin(\Theta_{U,V})\|_F \leq \frac{\|H\|_F}{\delta}.$$

Similarly, we could prove that it also holds for the  $2 \rightarrow 2$  norm.

$$\|\sin(\Theta_{U,V})\|_{2 \rightarrow 2} \leq \frac{\sqrt{k}\|H\|_{2 \rightarrow 2}}{\delta}.$$

In fact, for any unitary invariant norm  $\|\cdot\|$ , similar results hold.

## 7 A Useful Variant

One of the important applications of the Davis-Kahan Theorem is in statistics, where  $A$  captures the population information (thus is an unknown parameter, from the perspective of the frequentists), and  $H$  is the error due to sampling and randomness. In this setting, our goal is typically to show that



the sample version eigenvector will be quite close to the population one, thus being a good estimator. However, it is not natural to ask for eigengap between the population version and the sample version, since the sample version can actually take any value (although for some values, the probability is really small).

In light of this, there is a useful variant of Davis-Kahan Theorem, which has the same conclusion, but only requires the eigengap for one of the matrices,  $A$  or  $A + H$ , therefore allowing us to put restriction only on the population version.