

SVD Preliminaries

- $\mathbf{R}^n = \left\{ \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \mid x_i \text{ real} \right\}$

- inner product of $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$: $(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i y_i = \mathbf{x}^T \mathbf{y}$. Note: $\|\mathbf{x}\|_2 = \sqrt{(\mathbf{x}, \mathbf{x})}$.

- $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ orthogonal if $(\mathbf{x}, \mathbf{y}) = 0$.

- $\{\mathbf{q}_1, \dots, \mathbf{q}_k\} \subset \mathbf{R}^n$ orthonormal if $(\mathbf{q}_i, \mathbf{q}_j) = \delta_{i,j} = \begin{cases} 1, & i = j, \\ 0, & i \neq j. \end{cases}$

- A matrix $\mathbf{Q}_{n \times n}$ is orthogonal if $\mathbf{Q}^T = \mathbf{Q}^{-1}$.

- $\mathbf{Q}_{n \times n}$ is orthogonal \iff its columns $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ are an orthonormal set.

Proof:

$$(\mathbf{Q}^T \mathbf{Q})_{i,j} = (\text{row } i \text{ of } \mathbf{Q}^T, \text{ col } j \text{ of } \mathbf{Q}) = (\mathbf{q}_i, \mathbf{q}_j) = \delta_{i,j} \iff \{\mathbf{q}_1, \dots, \mathbf{q}_n\} \text{ orthonormal.}$$

- A set of orthonormal vectors $\{\mathbf{q}_1, \dots, \mathbf{q}_n\} \subset \mathbf{R}^n$ can be used as a basis for \mathbf{R}^n . In fact, it is particularly simple to represent an arbitrary vector $\mathbf{x} \in \mathbf{R}^n$ as a linear combination $\mathbf{x} = \sum_{i=1}^n c_i \mathbf{q}_i$,

for this is equivalent to $\mathbf{x} = \underbrace{[\mathbf{q}_1, \dots, \mathbf{q}_n]}_{\mathbf{Q}} \underbrace{\begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}}_{\mathbf{c}}$. Using the orthogonality of \mathbf{Q} , we have $\mathbf{c} = \mathbf{Q}^T \mathbf{x}$.

- \mathbf{Q} orthogonal $\implies \|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \mathbf{R}^n$.

Proof: $\|\mathbf{Q}\mathbf{x}\|_2 = \sqrt{(\mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{x})} = \sqrt{\mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x}} = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2$.

- $\mathbf{Q}_{n \times n}$ orthogonal $\implies \|\mathbf{A}\mathbf{Q}\|_2 = \|\mathbf{A}\|_2$ for arbitrary $\mathbf{A}_{m \times n}$, $\|\mathbf{Q}\mathbf{B}\|_2 = \|\mathbf{B}\|_2$ for arbitrary $\mathbf{B}_{n \times p}$. (These follow easily from the preceding result.)

- \mathbf{Q} orthogonal $\implies \kappa_2(\mathbf{Q}) = 1$. (Orthogonal matrices are perfectly conditioned in the 2-norm sense.)

- $\mathbf{A}_{n \times n}$ a real symmetric matrix $\implies \mathbf{A}$ has real eigenvalues $\lambda_1, \dots, \lambda_n$ and a corresponding set of orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$

$$\mathbf{A}\mathbf{q}_i = \lambda_i \mathbf{q}_i, \quad i = 1, \dots, n$$

$$(\mathbf{q}_i, \mathbf{q}_j) = \delta_{i,j}$$

Equivalently, in matrix form:

$$\mathbf{A} \underbrace{[\mathbf{q}_1, \dots, \mathbf{q}_n]}_{\mathbf{Q}} = \underbrace{[\mathbf{q}_1, \dots, \mathbf{q}_n]}_{\mathbf{Q}} \underbrace{\begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}}_{\mathbf{\Lambda}}$$

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T, \quad \mathbf{\Lambda} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}$$

i.e.,

$$\begin{aligned} \mathbf{A} [\mathbf{v}_1, \dots, \mathbf{v}_r | \mathbf{v}_{r+1}, \dots, \mathbf{v}_n] &= [\mathbf{u}_1, \dots, \mathbf{u}_r | \mathbf{u}_{r+1}, \dots, \mathbf{u}_m] \left[\begin{array}{ccc|ccc} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & & & & \\ & & & & & \\ & & & & \sigma_r & \\ \hline & & & 0 & & 0 \end{array} \right] \\ &= [\sigma_1 \mathbf{u}_1, \dots, \sigma_r \mathbf{u}_r | 0, \dots, 0], \end{aligned}$$

which is equivalent to

$$\begin{aligned} \mathbf{A} \mathbf{v}_1 &= \sigma_1 \mathbf{u}_1 \\ &\vdots \\ \mathbf{A} \mathbf{v}_r &= \sigma_r \mathbf{u}_r \\ \mathbf{A} \mathbf{v}_{r+1} &= 0 \\ &\vdots \\ \mathbf{A} \mathbf{v}_n &= 0 \end{aligned}$$

So what the SVD gives us is a pair of orthonormal bases

$$\begin{aligned} \{\mathbf{v}_1, \dots, \mathbf{v}_n\} &\text{ for } \mathbf{R}^n \\ \{\mathbf{u}_1, \dots, \mathbf{u}_m\} &\text{ for } \mathbf{R}^m \end{aligned}$$

in terms of which $\mathbf{A}_{m \times n}$ has a particularly simple (*decoupled*) description.

- $\mathbf{A} \mathbf{v}_i = \sigma_i \mathbf{u}_i, i = 1, \dots, r$
- $\text{range}(\mathbf{A}) \equiv \{\mathbf{y} \mid \mathbf{y} = \mathbf{A} \mathbf{x} \text{ for some } \mathbf{x}\} = \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$
 $\text{rank}(\mathbf{A}) \equiv \dim(\text{range}(\mathbf{A})) = r$
- $\text{null}(\mathbf{A}) \equiv \{\mathbf{x} \mid \mathbf{A} \mathbf{x} = 0\} = \text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ ($= \emptyset$ if $r = n$)
 $\dim(\text{null}(\mathbf{A})) = n - r$

Additional facts:

- $\mathbf{A}^T = \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T$ (\Rightarrow the singular values of \mathbf{A}^T are the same as those of \mathbf{A} and the left/right singular vectors of \mathbf{A}^T are the right/left singular vectors of \mathbf{A})
- $\mathbf{A}^T \mathbf{A} = \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma}) \mathbf{V}^T$ (\Rightarrow the eigenvalues and a corresponding set of orthonormal eigenvectors for $\mathbf{A}^T \mathbf{A}$ are $\{\sigma_1^2, \dots, \sigma_r^2, \underbrace{0, \dots, 0}_{n-r}\}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$)
- $\mathbf{A} \mathbf{A}^T = \mathbf{U} (\mathbf{\Sigma} \mathbf{\Sigma}^T) \mathbf{U}^T$ (\Rightarrow the eigenvalues and a corresponding set of orthonormal eigenvectors for $\mathbf{A} \mathbf{A}^T$ are $\{\sigma_1^2, \dots, \sigma_r^2, \underbrace{0, \dots, 0}_{m-r}\}$ and $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$).
- The SVD “layers” \mathbf{R}^n , the domain space of \mathbf{A} , according to sensitivity to \mathbf{A} , as described below:

$$\max_{\mathbf{x} \in \mathbf{R}^n, \|\mathbf{x}\|_2=1} \|\mathbf{A} \mathbf{x}\|_2 = \|\mathbf{A} \mathbf{v}_1\|_2 = \|\sigma_1 \mathbf{u}_1\|_2 = \sigma_1$$

$$\begin{aligned} \max_{\mathbf{x} \in \mathbf{R}^n, (\mathbf{x}, \mathbf{v}_1)=0, \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 &= \|\mathbf{A}\mathbf{v}_2\|_2 = \|\sigma_2 \mathbf{u}_2\|_2 = \sigma_2 \\ &\vdots \\ \max_{\mathbf{x} \in \mathbf{R}^n, (\mathbf{x}, \mathbf{v}_1)=\dots=(\mathbf{x}, \mathbf{v}_{i-1})=0, \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 &= \|\mathbf{A}\mathbf{v}_i\|_2 = \|\sigma_i \mathbf{u}_i\|_2 = \sigma_i \end{aligned}$$

Thus the dominant “mode” or effect of A is $\mathbf{v}_1 \rightarrow \sigma_1 \mathbf{u}_1$, followed by $\mathbf{v}_2 \rightarrow \sigma_2 \mathbf{u}_2$, etc.

- $A_{n \times n}$ is nonsingular if and only if $\sigma_n > 0$; $\kappa_2(\mathbf{A}_{n \times n}) = \sigma_1/\sigma_n$.

Applications of SVD

1. Solution of ill-conditioned system $\mathbf{A}_{n \times n} \mathbf{x} = \mathbf{b}$

$$m = n \Rightarrow \mathbf{\Sigma} = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{pmatrix}_{n \times n}$$

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \Rightarrow \mathbf{\Sigma}\mathbf{V}^T \mathbf{x} = \mathbf{U}^T \mathbf{b}$$

Defining $\xi = \mathbf{V}^T \mathbf{x}$, $\beta = \mathbf{U}^T \mathbf{b}$, we thus obtain an equivalent system with a *diagonal* coefficient matrix:

$$\mathbf{\Sigma} \xi = \beta. \tag{2}$$

Note that since \mathbf{V} and \mathbf{U} are orthogonal, $\mathbf{x} = \mathbf{V}\xi = \sum_{i=1}^n \xi_i \mathbf{v}_i$ and $\mathbf{b} = \mathbf{U}\beta = \sum_{i=1}^m \beta_i \mathbf{u}_i$. (Thus we’re expressing x and b in terms of our orthonormal bases for \mathbf{R}^n and \mathbf{R}^m , respectively.)

Solution of (2):

$$\xi_i = \frac{\beta_i}{\sigma_i}, \quad i = 1, \dots, n \quad (\Rightarrow \mathbf{x} = \sum_{i=1}^n \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i)$$

With SVD, we can now give a more complete answer to the stability question: If $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$ where $\|\delta\mathbf{b}\|_2 = \epsilon$, how large is $\delta\mathbf{x}$?

The corresponding transformed systems are $\mathbf{\Sigma}\xi = \beta$, $\mathbf{\Sigma}\delta\xi = \delta\beta$ where $\|\delta\beta\|_2 = \epsilon$. Thus

$$|\delta\xi_i| = \frac{|(\delta\beta)_i|}{\sigma_i} \leq \frac{\epsilon}{\sigma_i}.$$

Hence if σ_n is very small in comparison to the other σ_i ’s, we expect a large $\delta\xi_n$, in which case the error in $\mathbf{x} + \delta\mathbf{x}$ will be concentrated in the direction of \mathbf{v}_n ; the remaining portion of $\mathbf{x} + \delta\mathbf{x}$ may be quite accurate.

2. Least squares solution of an overdetermined linear system

Problem: Choose \mathbf{x} to minimize $Q(\mathbf{x}) = \|\mathbf{A}_{m \times n} \mathbf{x} - \mathbf{b}\|_2^2$, $m > n$

Solution via normal equations: $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$

(Denoting solution of normal equations by $\bar{\mathbf{x}}$, we have $Q(\bar{\mathbf{x}} + \delta \mathbf{x}) = Q(\bar{\mathbf{x}}) + \|\mathbf{A} \delta \mathbf{x}\|_2^2 \geq Q(\bar{\mathbf{x}})$.)

Solution via SVD (more stable)... Since $m > n$, Σ has the following configuration:

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_n & & \\ \hline & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix}_{m \times n}$$

Since \mathbf{U} is orthogonal,

$$\begin{aligned} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 &= \|\mathbf{U} \Sigma \mathbf{V}^T \mathbf{x} - \mathbf{b}\|_2 \\ &= \|\mathbf{U} (\Sigma \mathbf{V}^T \mathbf{x} - \mathbf{U}^T \mathbf{b})\|_2 \\ &= \|\Sigma \mathbf{V}^T \mathbf{x} - \mathbf{U}^T \mathbf{b}\|_2 \\ &= \|\Sigma \xi - \beta\|_2 \end{aligned}$$

where $\xi = \mathbf{V}^T \mathbf{x}$, $\beta = \mathbf{U}^T \mathbf{b}$. Now

$$\Sigma \xi - \beta = \begin{bmatrix} \sigma_1 \xi_1 - \beta_1 \\ \vdots \\ \sigma_n \xi_n - \beta_n \\ -\beta_{n+1} \\ \vdots \\ -\beta_m \end{bmatrix}$$

Thus

$$\min_{\mathbf{x} \in \mathbf{R}^n} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 = \sqrt{\sum_{i=n+1}^m \beta_i^2} = \sqrt{\sum_{i=n+1}^m (\mathbf{u}_i^T \mathbf{b})^2},$$

which is achieved for

$$\xi = \begin{bmatrix} \beta_1 / \sigma_1 \\ \vdots \\ \beta_n / \sigma_n \end{bmatrix}; \quad \text{equivalently, } \mathbf{x} = \sum_{i=1}^n \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i.$$

3. Lower rank approximation to $\mathbf{A}_{m \times n}$

Define

$$\mathbf{A}_k = \mathbf{U} \begin{bmatrix} \Sigma_k & | & 0 \\ \hline 0 & | & 0 \end{bmatrix} \mathbf{V}^T$$

where

$$\Sigma_k = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \\ & & & \ddots \\ & & & & 0 \end{bmatrix}.$$

Note that \mathbf{A}_k is a rank k approximation to \mathbf{A} (assuming $k < r$), and it should have good accuracy if the omitted singular values of \mathbf{A} (i.e., $\sigma_{k+1}, \dots, \sigma_r$) are small.

Observe that \mathbf{A}_k has an alternative representation:

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

where

$$\begin{aligned} \mathbf{U}_k &= [\mathbf{u}_1, \dots, \mathbf{u}_k]_{m \times k} \\ \mathbf{V}_k &= [\mathbf{v}_1, \dots, \mathbf{v}_k]_{n \times k} \end{aligned}$$

and that the total storage requirement for $\mathbf{U}_k, \mathbf{\Sigma}_k, \mathbf{V}_k$ is $(m + n + 1)k$ vs. mn for \mathbf{A} (or \mathbf{A}_k in non-factored form). If $k \ll \min\{m, n\}$ this can be an important saving. One application is image compression/coarsening (e.g., A = color-coded array of pixels).

Accuracy of \mathbf{A}_k in 2-norm and Frobenius norm:

$$\begin{aligned} \|A - A_k\|_2 &= \left\| U \left(\mathbf{\Sigma} - \left[\begin{array}{c|c} \mathbf{\Sigma}_k & 0 \\ \hline 0 & 0 \end{array} \right] \right) \mathbf{V}^T \right\|_2 \\ &= \left\| \mathbf{\Sigma} - \left[\begin{array}{c|c} \mathbf{\Sigma}_k & 0 \\ \hline 0 & 0 \end{array} \right] \right\|_2 \\ &= \sigma_{k+1} \end{aligned}$$

$$\begin{aligned} \|A - A_k\|_F &= \left\| U \left(\mathbf{\Sigma} - \left[\begin{array}{c|c} \mathbf{\Sigma}_k & 0 \\ \hline 0 & 0 \end{array} \right] \right) \mathbf{V}^T \right\|_F \\ &= \left\| \mathbf{\Sigma} - \left[\begin{array}{c|c} \mathbf{\Sigma}_k & 0 \\ \hline 0 & 0 \end{array} \right] \right\|_F \\ &= \sqrt{\sigma_{k+1}^2 + \dots + \sigma_r^2} \end{aligned}$$

Claim: A_k minimizes both $\|A - B\|_2, \|A - B\|_F$ over all rank k matrices B . (See Golub & Van Loan for proof.)

4. Orthogonal regression (also known as principle component analysis)

Suppose we have m data points $\{\mathbf{x}_i\}_{i=1}^m$ in \mathbf{R}^n , $m > n$, and we wish to fit them by a k -dimensional hyperplane \mathbf{H}_k :

$$\mathbf{H}_k = \mathbf{x}^* + \mathbf{S}_k.$$

Here \mathbf{x}^* is a point in \mathbf{R}^n and \mathbf{S}_k denotes a k -dimensional subspace of \mathbf{R}^n . The goal is to choose \mathbf{x}^* and \mathbf{S}_k to minimize $\sum_{i=1}^m d_i^2$ where d_i is the orthogonal distance from \mathbf{x}_i to \mathbf{H}_k .

Example: $n = 2, k = 1$. In this case, \mathbf{H}_k becomes

$$\mathbf{H}_1 = \left\{ \begin{pmatrix} x_1^* \\ x_2^* \end{pmatrix} + c \mathbf{q}_1 \right\}$$

where \mathbf{q}_1 is a unit vector in \mathbf{R}^2 . Thus we are fitting m points in 2-space by a straight line, with distances measured orthogonally (as opposed to vertically as in the case of least squares approximation).

Solution, in general:

$$\mathbf{x}^* = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \text{ (the mean of the data points).}$$

To determine S_k , form $A_{m \times n}$ with rows representing the data points \mathbf{x}_i , and compute its SVD. Then $\mathbf{S}_k = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$.

5. Web search

Suppose, as result of a keyword search, we have n potentially relevant web pages P_1, \dots, P_n . Which ones to report back to user? How to rank them? Useful technique: analyze link structure via SVD...

$$\begin{aligned} G &= (V, E) && \text{(graph representing links)} \\ V &= \{1, \dots, n\} && \text{(vertex } i \text{ represents } P_i) \\ E &= \{\overline{ij} \mid P_i \text{ points to } P_j\} && \text{(edges)} \end{aligned}$$

The corresponding adjacency matrix, denoted by $M_{n \times n}$, has as its i, j^{th} element

$$m_{i,j} = \begin{cases} 1, & \text{if } \overline{ij} \in E, \\ 0, & \text{if } \overline{ij} \notin E. \end{cases}$$

For any given subset of vertices $S \subset V$, there is a corresponding set of links T to other vertices. We describe sets S, T in terms of vectors $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ defined by

$$\begin{aligned} x_i &= \begin{cases} 1, & i \in S, \\ 0, & \text{otherwise,} \end{cases} \\ y_j &= \text{number of edges from } S \text{ to vertex } j. \end{aligned}$$

Note that $y_j = \sum_{i \in S} a_{i,j} = \sum_{i=1}^n x_i a_{i,j}$. Thus $\mathbf{y}^T = \mathbf{x}^T \mathbf{A}$; equivalently, $\mathbf{y} = \mathbf{A}^T \mathbf{x}$. This suggests the possibility of using an SVD of \mathbf{A}^T in order to summarize the primary information content of G in a condensed form.

Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ be the SVD of \mathbf{A} , in which case $\mathbf{A}^T = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T$ is the SVD of \mathbf{A}^T . The dominant modes in the transformation $\mathbf{A}^T \mathbf{x} = \mathbf{y}$ are then given by:

$$\begin{aligned} \mathbf{A}^T \mathbf{u}_1 &= \sigma_1 \mathbf{v}_1 \\ \mathbf{A}^T \mathbf{u}_2 &= \sigma_2 \mathbf{v}_2 \\ &\dots \end{aligned}$$

What we are looking for is a rapid decay in the σ_i 's so all but the first few are negligible.. The entries of $\mathbf{v}_1, \mathbf{v}_2, \dots$ yield "authority weights" for web pages P_1, P_2, \dots with respect to the given keyword. The vectors $\mathbf{u}_1, \mathbf{u}_2, \dots$ furnish a corresponding set of "hub weights".