

LINEAR ALGEBRA (XXVIII)

Singular Value Decompositions and Approximation

Again we will see many tools, in particular orthogonality, are of great use.

Review

Actually SVD can be viewed as an extension of the **spectral decompositions** of symmetric matrices to arbitrary matrices, both for the statement and for the method of the proof.

Theorem (spectral decomposition)

Let S be a **symmetric** $n \times n$ matrix with eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbb{R}$. Then there are n orthonormal vectors $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ such that

$$S = \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top + \lambda_2 \mathbf{v}_2 \mathbf{v}_2^\top + \cdots + \lambda_n \mathbf{v}_n \mathbf{v}_n^\top.$$

More precisely, if $\text{rank}(S) = r$, then $\lambda_{r+1} = \cdots = \lambda_n = 0$ and hence

$$S = \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top + \lambda_2 \mathbf{v}_2 \mathbf{v}_2^\top + \cdots + \lambda_r \mathbf{v}_r \mathbf{v}_r^\top.$$

Theorem (singular value decomposition)

Let A be an $m \times n$ matrix with **singular values** $\sigma_1, \dots, \sigma_r \in \mathbb{R}$ with $r = \text{rank}(A)$. Then there are orthonormal vectors $\mathbf{u}_1, \dots, \mathbf{u}_r \in \mathbb{R}^m$ and orthonormal $\mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{R}^n$ such that

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top.$$

Theorem (spectral decomposition)

Let S be a symmetric $n \times n$ matrix with eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbb{R}$. Then there is an orthogonal $n \times n$ matrix Q such that

$$S = Q \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} Q^T.$$

Theorem (singular value decomposition)

Let A be a real $m \times n$ matrix with $r = \text{rank}(A)$. Define an $m \times n$ matrix

$$\Sigma = \begin{bmatrix} D & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{bmatrix} \quad \text{for } D = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix},$$

where $\sigma_1, \dots, \sigma_r$ are the singular values of A . Then for some orthogonal $m \times m$ matrix U and orthogonal $n \times n$ matrix V

$$A = U\Sigma V^T.$$

Proof of SVD

Lemma

Let S be a real symmetric matrix. Then S has exactly $\text{rank}(S)$ nonzero eigenvalues (counting repetition).

Lemma

Let A be an $m \times n$ matrix. Then

$$\text{rank}(A^T A) = \text{rank}(A A^T) = \text{rank}(A).$$

Eigenvalues of $A^T A$ and AA^T

Lemma

*All eigenvalues of $A^T A$ and AA^T are **nonnegative**.*

*In otherwords, both $A^T A$ and AA^T are **positive semidefinite**.*

Eigenvalues of $A^T A$ and AA^T (cont'd)

Let

$$r = \text{rank}(A),$$

and by the previous lemmas:

- $\text{rank}(A^T A) = \text{rank}(AA^T) = r$, thus each of them has r nonzero eigenvalues.
- Every eigenvalue of $A^T A$ is nonnegative, so is every eigenvalue of AA^T .

Lemma

Both $A^T A$ and AA^T have the *same r positive eigenvalues (counting repetition)*. More precisely,

matrix	positive eigenvalues	eigenvectors
$A^T A$	$\lambda_1, \dots, \lambda_r$	$\mathbf{v}_1, \dots, \mathbf{v}_r$
AA^T	$\lambda_1, \dots, \lambda_r$	$\frac{1}{\sqrt{\lambda_1}} A\mathbf{v}_1, \dots, \frac{1}{\sqrt{\lambda_r}} A\mathbf{v}_r$

Definition

Let $\lambda_1, \dots, \lambda_r$ be the nonzero (hence positive) eigenvalues of $A^T A$ (hence the nonzero eigenvalues of AA^T). For every $i \in [r]$, we define

$$\sigma_i = \sqrt{\lambda_i}.$$

Then $\sigma_1, \dots, \sigma_r$ are the singular values of A .

Theorem

Let A be a real $m \times n$ matrix with $r = \text{rank}(A)$. Define an $m \times n$ matrix

$$\Sigma = \begin{bmatrix} D & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{bmatrix} \quad \text{for } D = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix},$$

where $\sigma_1, \dots, \sigma_r$ are the singular values of A . Then for some orthogonal $m \times m$ matrix U and orthogonal $n \times n$ matrix V

$$A = U\Sigma V^T.$$

Let

$$\lambda_1, \dots, \lambda_r, \lambda_{r+1}, \dots, \lambda_n$$

be the eigenvalues of $A^T A$ with

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$$

and $\lambda_{r+1} = \dots = \lambda_n = 0$. Moreover, let

$$\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$$

be the corresponding orthonormal eigenvectors. It follows that

$$V = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n].$$

is an orthogonal $n \times n$ matrix.

Next for every $i \in [r]$ let

$$\mathbf{u}_i = \frac{1}{\sqrt{\lambda_i}} A \mathbf{v}_i \in \mathbb{R}^m.$$

Thus

$$\mathbf{u}_1, \dots, \mathbf{u}_r$$

are orthonormal eigenvectors of AA^T with corresponding eigenvalues $\lambda_1, \dots, \lambda_r$. We extend $\mathbf{u}_1, \dots, \mathbf{u}_r$ to a basis

$$\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_m$$

for \mathbb{R}^m . By Gram-Schmidt, we can further assume that they are orthonormal. Therefore,

$$U = [\mathbf{u}_1 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m]$$

is orthogonal.

Let $i \in [n]$. If $i \leq r$, then

$$A\mathbf{v}_i = \sqrt{\lambda_i}\mathbf{u}_i.$$

Otherwise, $i > r$, and then $\lambda_i = 0$. Since \mathbf{v}_i is an eigenvector of $A^\top A$ for λ_i ,

$$(A^\top A)\mathbf{v}_i = \lambda_i\mathbf{v}_i = \mathbf{0}.$$

It follows that

$$0 = \mathbf{v}_i^\top (A^\top A)\mathbf{v}_i = A\mathbf{v}_i \cdot A\mathbf{v}_i = \|A\mathbf{v}_i\|,$$

thereby $A\mathbf{v}_i = \mathbf{0}$. We can thus deduce

$$\begin{aligned} AV &= [A\mathbf{v}_1 \ \dots \ A\mathbf{v}_r \ A\mathbf{v}_{r+1} \ \dots \ A\mathbf{v}_n] = [\sqrt{\lambda_1}\mathbf{u}_1 \ \dots \ \sqrt{\lambda_r}\mathbf{u}_r \ \mathbf{0} \ \dots \ \mathbf{0}] \\ &= [\mathbf{u}_1 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sqrt{\lambda_1} & & & & & \\ & \ddots & & & & \\ & & \sqrt{\lambda_r} & & & \\ & & & 0 & & \\ & & & & \ddots & \\ & & & & & 0 \end{bmatrix} = U\Sigma. \end{aligned}$$

Since V is orthogonal, i.e., $V^{-1} = V^T$, we conclude

$$A = U\Sigma V^T.$$



Bases and matrices in the SVD

Singular value decomposition of an arbitrary A

Theorem

Let A be an $m \times n$ matrix with $r = \text{rank}(A)$. Define an $m \times n$ matrix

$$\Sigma = \begin{bmatrix} D & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{bmatrix} \quad \text{for } D = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix},$$

where $\sigma_1, \dots, \sigma_r$ are the singular values of A . Then for some orthogonal $m \times m$ matrix U and orthogonal $n \times n$ matrix V

$$A = U\Sigma V^T.$$

In particular, there are $\mathbf{u}_1, \dots, \mathbf{u}_r \in \mathbb{R}^m$ and $\mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{R}^n$ such that

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T.$$

From

$$\begin{aligned}
 A &= [\mathbf{u}_1 \ \cdots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \cdots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_r^\top \\ \mathbf{v}_{r+1}^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix} \\
 &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top
 \end{aligned}$$

we deduce:

Corollary

vectors	orthonormal basis for the vector space
$\mathbf{u}_1, \dots, \mathbf{u}_r$	$C(A)$
$\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$	$N(A^\top)$
$\mathbf{v}_1, \dots, \mathbf{v}_r$	$C(A^\top)$
$\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$	$N(A)$

Proof of $\mathbf{u}_1, \dots, \mathbf{u}_r$ as a basis for $C(A)$

As $\dim(C(A)) = \text{rank}(A)$, we only need to show

$$C(A) \subseteq \text{span}(\{\mathbf{u}_1, \dots, \mathbf{u}_r\}).$$

Recall

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T.$$

So (can you see?) every column of A is a linear combination of $\mathbf{u}_1, \dots, \mathbf{u}_r$. \square

Proof of $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$ as a basis for $\mathbf{N}(A)$

As $\dim(\mathbf{N}(A)) = n - \text{rank}(A) = n - r$, we only need to show

$$A\mathbf{v}_i = \mathbf{0}$$

for every $r + 1 \leq i \leq n$. Again by

$$A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top,$$

together with the orthogonality of $\mathbf{v}_1, \dots, \mathbf{v}_r, \dots, \mathbf{v}_i, \dots, \mathbf{v}_n$, we obtain

$$\begin{aligned} A\mathbf{v}_i &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top \mathbf{v}_i + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top \mathbf{v}_i + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top \mathbf{v}_i \\ &= \sigma_1 \mathbf{u}_1 (\mathbf{v}_1 \cdot \mathbf{v}_i) + \sigma_2 \mathbf{u}_2 (\mathbf{v}_2 \cdot \mathbf{v}_i) + \cdots + \sigma_r \mathbf{u}_r (\mathbf{v}_r \cdot \mathbf{v}_i) \\ &= \mathbf{0}. \end{aligned}$$



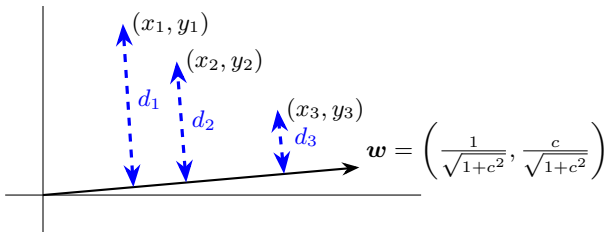
Geometry and Applications of SVD

Fitting a set of data by a line

Assume that we have m points in \mathbb{R}^2

$$(x_1, y_1), \dots, (x_m, y_m).$$

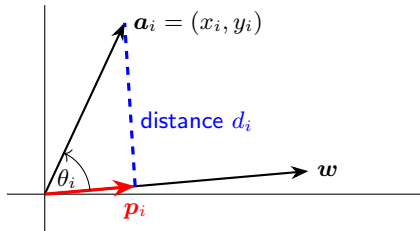
We want to find a line $y = cx$ for some $c \in \mathbb{R}$ such that (x_i, y_i) s are close to it as much as possible. Equivalently, this line can be determined by a unit vector $w \in \mathbb{R}^2$



such that $d_1^2 + d_2^2 + \dots + d_m^2$ is minimum.

Fitting a set of data by a line (cont'd)

For $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^2$ we want to determine a **unit vector** $\mathbf{w} \in \mathbb{R}^2$



to **minimize**

$$\sum_{i \in [m]} d_i^2 = \sum_{i \in [m]} (\|\mathbf{a}_i\|^2 - \|\mathbf{p}_i\|^2) = \sum_{i \in [m]} \|\mathbf{a}_i\|^2 - \sum_{i \in [m]} \|\mathbf{p}_i\|^2.$$

Equivalently, we **maximize**

$$\sum_{i \in [m]} \|\mathbf{p}_i\|^2 = \sum_{i \in [m]} (\|\mathbf{a}_i\| \cos \theta_i)^2 = \sum_{i \in [m]} \left(\|\mathbf{a}_i\| \frac{\mathbf{w} \cdot \mathbf{a}_i}{\|\mathbf{w}\| \|\mathbf{a}_i\|} \right)^2 = \sum_{i \in [m]} (\mathbf{w} \cdot \mathbf{a}_i)^2.$$

Each data now has n features. Given

$$\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$$

we want to find a unit vector $\mathbf{w} \in \mathbb{R}^n$ to minimize

$$\sum_{i \in [m]} (\text{the distance between } \mathbf{a}_i \text{ and the line along } \mathbf{w})^2.$$

It is equivalent to maximizing

$$\sum_{i \in [m]} (\mathbf{w} \cdot \mathbf{a}_i)^2,$$

whose precise proof is standard by now and will come later.

Define an $m \times n$ matrix by

$$A = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}.$$

We look for a unit vector $\mathbf{w} \in \mathbb{R}^n$ such that

$$\sum_{i \in [m]} (\mathbf{w} \cdot \mathbf{a}_i)^2 = \sum_{i \in [m]} (\mathbf{a}_i \cdot \mathbf{w})^2 = \sum_{i \in [m]} \left(\mathbf{a}_i^\top \mathbf{w} \right)^2 = \|A\mathbf{w}\|^2$$

is maximized.

First singular vector maximizes $\|Av\|$

Theorem

Let A be an $m \times n$ matrix and σ_1 be its *largest* singular value. Then

$$\max_{\substack{\mathbf{w} \in \mathbb{R}^n \text{ with} \\ \|\mathbf{w}\| = 1}} \|A\mathbf{w}\|^2 = \sigma_1^2.$$

Moreover, the maximum is attained when $\mathbf{w} = \mathbf{v}_1$ where \mathbf{v}_1 appears in the SVD of A

$$A = U\Sigma V = U\Sigma \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix}.$$

Orthogonality, revisited

Orthonormal basis and coordinate vectors

Consider an orthonormal basis $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$. Then every $\mathbf{w} \in \mathbb{R}^n$ can be expressed as a linear combination

$$\mathbf{w} = c_1 \mathbf{v}_1 + \dots + c_n \mathbf{v}_n,$$

i.e., the **coordinate vector** of \mathbf{w} with respect to $\mathbf{v}_1, \dots, \mathbf{v}_n$ is $(c_1, \dots, c_n) \in \mathbb{R}^n$.

Using the orthonormality of \mathbf{v}_i s, each $c_i = \mathbf{v}_i \cdot \mathbf{w}$ and hence

$$c_i \mathbf{v}_i = (\mathbf{v}_i \cdot \mathbf{w}) \mathbf{v}_i = \mathbf{v}_i^T \mathbf{w} \mathbf{v}_i = \mathbf{v}_i \mathbf{v}_i^T \mathbf{w},$$

i.e.,

$$\begin{aligned} \mathbf{w} &= (\mathbf{v}_1 \cdot \mathbf{w}) \mathbf{v}_1 + \dots + (\mathbf{v}_n \cdot \mathbf{w}) \mathbf{v}_n = \mathbf{v}_1 \mathbf{v}_1^T \mathbf{w} + \dots + \mathbf{v}_n \mathbf{v}_n^T \mathbf{w} \\ &= [\mathbf{v}_1 \ \dots \ \mathbf{v}_n] \begin{bmatrix} \mathbf{v}_1^T \mathbf{w} \\ \vdots \\ \mathbf{v}_n^T \mathbf{w} \end{bmatrix} = \underbrace{[\mathbf{v}_1 \ \dots \ \mathbf{v}_n]}_Q \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} \mathbf{w} = Q Q^T \mathbf{w}. \end{aligned}$$

This is precisely that

the least squares solution of $Q\mathbf{x} = \mathbf{b}$ is $\hat{\mathbf{x}} = Q^T \mathbf{b}$.

Let $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^n$ be orthonormal and

$$Q = [\mathbf{q}_1 \ \dots \ \mathbf{q}_n].$$

This is equivalent to that Q is an orthogonal matrix.

So the least squares solution of $Q\mathbf{x} = \mathbf{b}$ is $\hat{\mathbf{x}} = Q^\top \mathbf{b}$. The projection matrix is $QQ^\top = I$, hence

$$\mathbf{b} = QQ^\top \mathbf{b} = [\mathbf{q}_1 \ \dots \ \mathbf{q}_n] \begin{bmatrix} \mathbf{q}_1^\top \mathbf{b} \\ \vdots \\ \mathbf{q}_n^\top \mathbf{b} \end{bmatrix} = \mathbf{q}_1 \mathbf{q}_1^\top \mathbf{b} + \dots + \mathbf{q}_n \mathbf{q}_n^\top \mathbf{b}.$$

\mathbf{b} is the sum of projections of \mathbf{b} onto every line $\text{span}(\{\mathbf{q}_i\})$.

Second look through change of basis

We consider the vector space $V = \mathbb{R}^n$ and two of its bases

$$e_1, \dots, e_n \quad \text{and} \quad v_1, \dots, v_n.$$

So the change of basis matrix M from e_1, \dots, e_n to v_1, \dots, v_n satisfies

$$[v_1 \ \dots \ v_n] = [e_1 \ \dots \ e_n]M.$$

It is easy to see (why?)

$$M = [v_1 \ \dots \ v_n].$$

Second look through change of basis (cont'd)

Consider an arbitrary $\mathbf{w} = (w_1, \dots, w_n) \in \mathbf{V}$. So the coordinate vector of \mathbf{w} with respect to $\mathbf{e}_1, \dots, \mathbf{e}_n$ is $(w_1, \dots, w_n) \in \mathbb{R}^n$, precisely \mathbf{w} itself. Recall:

Theorem

Let $\bar{\mathbf{v}}$ and $\bar{\mathbf{v}}'$ be two bases for \mathbf{V} . Then for every $\mathbf{v} \in \mathbb{R}^n$

$$T_{\bar{\mathbf{v}}}(\mathbf{v}) = MT_{\bar{\mathbf{v}}'}(\mathbf{v}).$$

Hence the coordinate vector of \mathbf{w} with respect to $\mathbf{v}_1, \dots, \mathbf{v}_n$ is

$$M^{-1}\mathbf{w}.$$

If $\mathbf{v}_1, \dots, \mathbf{v}_n$ are orthonormal, i.e.,

$$M = Q = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n]$$

is orthogonal, then

$$M^{-1}\mathbf{w} = Q^T\mathbf{w} \quad \text{and} \quad \mathbf{w} = QQ^T\mathbf{w} = \mathbf{v}_1\mathbf{v}_1^T\mathbf{w} + \dots + \mathbf{v}_n\mathbf{v}_n^T\mathbf{w}.$$

Recall

Theorem

Let Q be an $n \times n$ orthogonal matrix and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

- (i) $\|Q\mathbf{x}\| = \|\mathbf{x}\|$, hence $\|Q^T\mathbf{x}\| = \|\mathbf{x}\|$.
- (ii) $Q\mathbf{x} \cdot Q\mathbf{y} = \mathbf{x} \cdot \mathbf{y}$, hence $Q^T\mathbf{x} \cdot Q^T\mathbf{y} = \mathbf{x} \cdot \mathbf{y}$.

Let $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ be an orthonormal basis, or equivalently

$$Q = [\mathbf{v}_1 \ \dots \ \mathbf{v}_n]$$

is orthogonal. Consider an arbitrary vector $\mathbf{w} \in \mathbb{R}^n$.

- (i) The coordinate vector of \mathbf{w} with respect to $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ is $Q^T \mathbf{w}$, i.e.,

$$(\mathbf{v}_1^T \mathbf{w}, \dots, \mathbf{v}_n^T \mathbf{w}).$$

- (ii) The length of the coordinate vector of \mathbf{w} with respect to $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ equals to the length of \mathbf{w} , i.e.,

$$\|\mathbf{w}\| = \|Q^T \mathbf{w}\| = \sqrt{(\mathbf{v}_1^T \mathbf{w})^2 + \dots + (\mathbf{v}_n^T \mathbf{w})^2}.$$

Theorem

Let A be an $m \times n$ matrix and σ_1 be its *largest* singular value. Then

$$\max_{\substack{\mathbf{w} \in \mathbb{R}^n \text{ with} \\ \|\mathbf{w}\| = 1}} \|A\mathbf{w}\| = \sigma_1.$$

Moreover, the maximum is attained when $\mathbf{w} = \mathbf{v}_1$, where \mathbf{v}_1 appears in the SVD of A

$$A = U\Sigma V = U\Sigma \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix}.$$

Let $\lambda_1 > \dots > \lambda_r$ be the positive eigenvalues of $A^\top A$. For every $i \in [r]$, we define $\sigma_i = \sqrt{\lambda_i}$. Then $\sigma_1 > \dots > \sigma_r$ are the singular values of A . Moreover, for some orthonormal $\mathbf{u}_1, \dots, \mathbf{u}_m \in \mathbb{R}^m$ and $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$

$$A = [\mathbf{u}_1 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_r & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_r^\top \\ \mathbf{v}_{r+1}^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix}$$

$$= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top$$

Let $\mathbf{w} \in \mathbb{R}^n$ be a unit vector. So $\mathbf{w} = c_1 \mathbf{v}_1 + \cdots + c_n \mathbf{v}_n$ for some $c_1, \dots, c_n \in \mathbb{R}$ with $\sum_{i \in [n]} c_i^2 = 1$. Hence

$$\begin{aligned} A\mathbf{w} &= (\sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top) (c_1 \mathbf{v}_1 + \cdots + c_n \mathbf{v}_n) \\ &= \sum_{i \in [r], j \in [n]} \sigma_i c_j \mathbf{u}_i \mathbf{v}_i^\top \mathbf{v}_j = \sum_{i \in [r]} \sigma_i c_i \mathbf{u}_i. \end{aligned}$$

Then the coordinate vector of $A\mathbf{w}$ with respect to the orthonormal basis $\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ is $(\sigma_1 c_1, \dots, \sigma_r c_r, 0, \dots, 0)$. We deduce

$$\begin{aligned} \|A\mathbf{w}\|^2 &= \|(\sigma_1 c_1, \dots, \sigma_r c_r, 0, \dots, 0)\|^2 = \sum_{i \in [r]} \sigma_i^2 c_i^2 \\ &\leq \sigma_1^2 \sum_{i \in [r]} c_i^2 \leq \sigma_1^2 \sum_{i \in [n]} c_i^2 = \sigma_1^2. \end{aligned}$$

And the equality can be attained by choosing $\mathbf{w} = \mathbf{v}_1$.



Fitting data by a k -dimension subspace

Fitting a set of data by a k -dimension subspace

Fix a constant $k \geq 1$. Given

$$\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$$

we want to find a subspace $\mathbf{W} \subseteq \mathbb{R}^n$ with $\dim(\mathbf{W}) = k$ to minimize

$$\sum_{i \in [m]} (\text{the distance between } \mathbf{a}_i \text{ and } \mathbf{W})^2.$$

Observe that \mathbf{W} can be determined by any linearly independent

$$\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbf{W}.$$

And by Gram-Schmidt, $\mathbf{w}_1, \dots, \mathbf{w}_k$ can be chosen as orthonormal.

The distance between \mathbf{a}_i and \mathbf{W} is

the distance between \mathbf{a}_i and the projection \mathbf{p}_i of \mathbf{a}_i onto \mathbf{W} .

1. We want to compute the projection \mathbf{p} of \mathbf{b} onto a subspace

$$\mathbf{V} = \text{span}(\{\mathbf{a}_1, \dots, \mathbf{a}_n\})$$

with linearly independent $\mathbf{a}_1, \dots, \mathbf{a}_n$. In particular $\mathbf{p} \in \mathbf{V}$.

2. We define \mathbf{p} as a vector whose error vector $\mathbf{e} = \mathbf{b} - \mathbf{p}$ is perpendicular to \mathbf{V} .
3. We prove that the projection matrix is $P = A(A^\top A)^{-1}A^\top$, i.e.,

$$\mathbf{p} = P\mathbf{b} = A(A^\top A)^{-1}A^\top \mathbf{b}.$$

Here, we have verified that $(A^\top A)^{-1}$ really exists provided $\text{rank}(A) = n$.

This also implies that \mathbf{p} is the **unique** vector whose error vector $\mathbf{e} = \mathbf{b} - \mathbf{p}$ is perpendicular to \mathbf{V} .

4. As an exercise, we can show that for every $\mathbf{v} \in \mathbf{V}$

$$\|\mathbf{b} - \mathbf{v}\| = \min_{\mathbf{u} \in \mathbf{V}} \|\mathbf{b} - \mathbf{u}\| \iff \mathbf{v} = \mathbf{p}.$$

We want to determine \mathbf{W} by choosing orthonormal

$$\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^n.$$

To that end, let

$$\mathbf{Q} = [\mathbf{w}_1 \ \dots \ \mathbf{w}_k].$$

So

$$\mathbf{Q}^\top \mathbf{Q} = \begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_k^\top \end{bmatrix} [\mathbf{w}_1 \ \dots \ \mathbf{w}_k] = \begin{bmatrix} \mathbf{w}_1^\top \mathbf{w}_1 & \cdots & \mathbf{w}_1^\top \mathbf{w}_k \\ \vdots & \ddots & \vdots \\ \mathbf{w}_k^\top \mathbf{w}_1 & \cdots & \mathbf{w}_k^\top \mathbf{w}_k \end{bmatrix} = \mathbf{I}_{k \times k},$$

and

$$\mathbf{Q} \mathbf{Q}^\top = [\mathbf{w}_1 \ \dots \ \mathbf{w}_k] \begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_k^\top \end{bmatrix} = \mathbf{w}_1 \mathbf{w}_1^\top + \cdots + \mathbf{w}_k \mathbf{w}_k^\top.$$

Then the projection of \mathbf{a}_i onto $\mathbf{W} = \text{span}(\{\mathbf{w}_1, \dots, \mathbf{w}_k\}) = \mathbf{C}(Q)$ is

$$\begin{aligned} \mathbf{p}_i &= Q(Q^\top Q)^{-1}Q^\top \mathbf{a}_i = QQ^\top \mathbf{a}_i = \mathbf{w}_1\mathbf{w}_1^\top \mathbf{a}_i + \dots + \mathbf{w}_k\mathbf{w}_k^\top \mathbf{a}_i \\ &= (\mathbf{w}_1 \cdot \mathbf{a}_i)\mathbf{w}_1 + \dots + (\mathbf{w}_k \cdot \mathbf{a}_i)\mathbf{w}_k = \sum_{j \in [k]} (\mathbf{w}_j \cdot \mathbf{a}_i)\mathbf{w}_j. \end{aligned}$$

This is exactly

the sum of projections of \mathbf{a}_i onto the lines $\text{span}(\{\mathbf{w}_1\}), \dots, \text{span}(\{\mathbf{w}_k\})$.

So the distance between \mathbf{a}_i and \mathbf{W} is $d_i = \|\mathbf{a}_i - \mathbf{p}_i\|$, where $\mathbf{a}_i - \mathbf{p}_i$ is the previous **error vector** \mathbf{e} . Then (basically a previous exercise)

$$\begin{aligned} & \|\mathbf{a}_i - \mathbf{p}_i\|^2 \\ &= (\mathbf{a}_i - \mathbf{p}_i) \cdot (\mathbf{a}_i - \mathbf{p}_i) = (\mathbf{a}_i - \mathbf{p}_i)^\top (\mathbf{a}_i - \mathbf{p}_i) \\ &= (\mathbf{a}_i^\top - \mathbf{p}_i^\top) (\mathbf{a}_i - \mathbf{p}_i) = \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{p}_i^\top \mathbf{a}_i - \mathbf{a}_i^\top \mathbf{p}_i + \mathbf{p}_i^\top \mathbf{p}_i \\ &= \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{a}_i^\top \mathbf{p}_i - \mathbf{a}_i^\top \mathbf{p}_i + \mathbf{p}_i^\top \mathbf{p}_i = \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{p}_i^\top \mathbf{p}_i - 2(\mathbf{a}_i^\top \mathbf{p}_i - \mathbf{p}_i^\top \mathbf{p}_i) \\ &= \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{p}_i^\top \mathbf{p}_i - 2(\mathbf{a}_i - \mathbf{p}_i)^\top \mathbf{p}_i = \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{p}_i^\top \mathbf{p}_i - 2(\mathbf{a}_i - \mathbf{p}_i) \cdot \mathbf{p}_i \\ &= \mathbf{a}_i^\top \mathbf{a}_i - \mathbf{p}_i^\top \mathbf{p}_i \quad (\text{by } (\mathbf{a}_i - \mathbf{p}_i) \perp \mathbf{p}_i) \\ &= \|\mathbf{a}_i\|^2 - \|\mathbf{p}_i\|^2. \end{aligned}$$

By

$$d_i^2 = \|\mathbf{a}_i - \mathbf{p}_i\|^2 = \|\mathbf{a}_i\|^2 - \|\mathbf{p}_i\|^2 \quad \text{and} \quad \mathbf{p}_i = \sum_{j \in [k]} (\mathbf{w}_j \cdot \mathbf{a}_i) \mathbf{w}_j,$$

it is equivalent to maximizing

$$\sum_{i \in [m]} \|\mathbf{p}_i\|^2 = \sum_{i \in [m]} \left\| \sum_{j \in [k]} (\mathbf{w}_j \cdot \mathbf{a}_i) \mathbf{w}_j \right\|^2 = \sum_{i \in [m]} \sum_{j \in [k]} (\mathbf{w}_j \cdot \mathbf{a}_i)^2.$$

When $k = 1$ and $\mathbf{w}_1 = \mathbf{w}$, this is what we claimed: minimizing

$$\sum_{i \in [m]} (\text{the distance between } \mathbf{a}_i \text{ and the line along a unit vector } \mathbf{w})^2.$$

is equivalent to maximizing

$$\sum_{i \in [m]} (\mathbf{w} \cdot \mathbf{a}_i)^2,$$

Recall we combine m points $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$ into an $m \times n$ matrix

$$A = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}.$$

We look for orthonormal vectors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^n$ such that

$$\sum_{i \in [m]} \sum_{j \in [k]} (\mathbf{w}_j \cdot \mathbf{a}_i)^2 = \sum_{j \in [k]} \sum_{i \in [m]} (\mathbf{w}_j \cdot \mathbf{a}_i)^2 = \sum_{j \in [k]} \|A\mathbf{w}_j\|^2$$

is maximized.

Recall:

$$\begin{aligned}
 A &= [\mathbf{u}_1 \ \cdots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \cdots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_r^\top \\ \mathbf{v}_{r+1}^\top \\ \vdots \\ \mathbf{v}_n^\top \end{bmatrix} \\
 &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top
 \end{aligned}$$

Theorem

Assume $k \leq r$. Then

$$\sum_{j \in [k]} \|A \mathbf{w}_j\|^2 \leq \sum_{j \in [k]} \sigma_j^2.$$

And the equality holds for $\mathbf{w}_1 = \mathbf{v}_1, \dots$, and $\mathbf{w}_k = \mathbf{v}_k$.