

Definition 1: Eigendecomposition & Diagonalization

Suppose an n by n matrix A has n linearly independent eigenvectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. Put them into the columns of the **eigenvector matrix** X . Then $\Lambda \equiv X^{-1}AX$ is the **eigenvalue matrix**. A has been **diagonalized**.

$$\Lambda \equiv X^{-1}AX \quad AX = X\Lambda \quad A = X\Lambda X^{-1}$$

$X\Lambda X^{-1}$ is called the **eigendecomposition** of A . But what does the eigenvalue matrix Λ look like?

Why do we say A has been "diagonalized"?

Lemma 1: Eigenvalue Matrix Is Diagonal With Eigenvalues In Diagonal

Suppose an n by n matrix A has n linearly independent eigenvectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. Then the eigenvalue matrix $\Lambda \equiv X^{-1}AX$ is a diagonal matrix with the eigenvalues in the diagonal.

Proof:

AX and $X\Lambda$ are both equal to $[\lambda_1\mathbf{x}_1 \ \dots \ \lambda_n\mathbf{x}_n]$:

$$\begin{aligned} AX &= \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1\mathbf{x}_1 & \mathbf{a}_1\mathbf{x}_2 & \dots & \mathbf{a}_1\mathbf{x}_n \\ \mathbf{a}_2\mathbf{x}_1 & \mathbf{a}_2\mathbf{x}_2 & \dots & \mathbf{a}_2\mathbf{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_n\mathbf{x}_1 & \mathbf{a}_n\mathbf{x}_2 & \dots & \mathbf{a}_n\mathbf{x}_n \end{bmatrix} \\ &= \begin{bmatrix} A\mathbf{x}_1 & \dots & A\mathbf{x}_n \end{bmatrix} = \begin{bmatrix} \lambda_1\mathbf{x}_1 & \dots & \lambda_n\mathbf{x}_n \end{bmatrix} \\ X\Lambda &= \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{2,1} & \dots & x_{n,1} \\ x_{1,2} & x_{2,2} & \dots & x_{n,2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1,n} & x_{2,n} & \dots & x_{n,n} \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} \\ &= \begin{bmatrix} x_{1,1}\lambda_1 & x_{2,1}\lambda_2 & \dots & x_{n,1}\lambda_n \\ x_{1,2}\lambda_1 & x_{2,2}\lambda_2 & \dots & x_{n,2}\lambda_n \\ \vdots & \vdots & \ddots & \vdots \\ x_{1,n}\lambda_1 & x_{2,n}\lambda_2 & \dots & x_{n,n}\lambda_n \end{bmatrix} = \begin{bmatrix} \lambda_1\mathbf{x}_1 & \dots & \lambda_n\mathbf{x}_n \end{bmatrix} \end{aligned}$$

The eigenvectors are linearly independent, so X is invertible and we have

$$X^{-1}AX = \Lambda \quad AX = X\Lambda \quad A = X\Lambda X^{-1}$$

■

Let's try an example:

$$A = \begin{bmatrix} 1 & 2 \\ 4 & 3 \end{bmatrix} \rightarrow 0 = \begin{vmatrix} 1-\lambda & 2 \\ 4 & 3-\lambda \end{vmatrix} = (1-\lambda)(3-\lambda) - 8 = \lambda^2 - 4\lambda + 3 - 8$$

$$= \lambda^2 - 4\lambda - 5 = (\lambda + 1)(\lambda - 5) \implies \lambda_1 = -1 \quad \lambda_2 = 5$$

Let's check: $|A| = 1 \cdot 3 - 2 \cdot 4 = -5 = \lambda_1 \cdot \lambda_2 \quad \text{tr}(A) = 4 = \lambda_1 + \lambda_2$

$$\lambda_1 = -1 \rightarrow \begin{bmatrix} 1-\lambda_1 & 2 \\ 4 & 3-\lambda_1 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 4 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \rightarrow \mathbf{x}_1 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$A\mathbf{x}_1 = \begin{bmatrix} 1 & 2 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1+2 \\ -4+3 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \lambda_1 \mathbf{x}_1$$

$$\lambda_2 = 5 \rightarrow \begin{bmatrix} 1-\lambda_2 & 2 \\ 4 & 3-\lambda_2 \end{bmatrix} = \begin{bmatrix} -4 & 2 \\ 4 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -\frac{1}{2} \\ 0 & 0 \end{bmatrix} \rightarrow \mathbf{x}_2 = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}$$

$$A\mathbf{x}_2 = \begin{bmatrix} 1 & 2 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}+2 \\ 2+3 \end{bmatrix} = \begin{bmatrix} \frac{5}{2} \\ 5 \end{bmatrix} = \lambda_2 \mathbf{x}_2$$

Let's check that our eigenvectors are independent:

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$\begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} -1 & \frac{1}{2} \\ 0 & \frac{3}{2} \end{bmatrix} \rightarrow \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

\implies No free variables and nullspace = $\{\mathbf{0}\}$

And let's calculate X^{-1} :

$$X^{-1} = \frac{1}{|X|} \text{adj}(X) = \frac{1}{-1-\frac{1}{2}} \begin{bmatrix} 1 & -\frac{1}{2} \\ -1 & -1 \end{bmatrix} = -\frac{2}{3} \begin{bmatrix} 1 & -\frac{1}{2} \\ -1 & -1 \end{bmatrix}$$

And check the inverse calculation:

$$X^{-1}X = -\frac{2}{3} \begin{bmatrix} 1 & -\frac{1}{2} \\ -1 & -1 \end{bmatrix} \begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} = -\frac{2}{3} \begin{bmatrix} -1-\frac{1}{2} & \frac{1}{2}-\frac{1}{2} \\ 1-1 & -\frac{1}{2}-1 \end{bmatrix} = -\frac{2}{3} \begin{bmatrix} -\frac{3}{2} & 0 \\ 0 & -\frac{3}{2} \end{bmatrix} = I$$

And, finally, let's diagonalize A :

$$X\Lambda X^{-1} = \begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \left(-\frac{2}{3}\right) \begin{bmatrix} 1 & -\frac{1}{2} \\ -1 & -1 \end{bmatrix}$$

$$= \left(-\frac{2}{3}\right) \begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & -\frac{1}{2} \\ -1 & -1 \end{bmatrix}$$

$$= \left(-\frac{2}{3}\right) \begin{bmatrix} -1 & \frac{1}{2} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & \frac{1}{2} \\ -5 & -5 \end{bmatrix} = \left(-\frac{2}{3}\right) \begin{bmatrix} 1-\frac{5}{2} & -\frac{1}{2}-\frac{5}{2} \\ -1-5 & \frac{1}{2}-5 \end{bmatrix}$$

$$= \left(-\frac{2}{3}\right) \begin{bmatrix} -\frac{3}{2} & -\frac{6}{2} \\ -6 & -\frac{9}{2} \end{bmatrix} = \frac{2}{3} \begin{bmatrix} \frac{3}{2} & \frac{6}{2} \\ 6 & \frac{9}{2} \end{bmatrix} = A$$

So we can't just diagonalize any ol' square matrix? How do we know which matrices have independent eigenvectors? Please don't tell me that we have to compute the eigenvectors to see if they're independent.

Lemma 2: Independent Eigenvectors From Distinct Eigenvalues

The eigenvectors from distinct eigenvalues are linearly independent.

Proof:

Suppose $\sum_{i=1}^n c_i \mathbf{x}_i = \mathbf{0}$. Multiply by A and λ_n to get two equations:

$$\mathbf{0} = A \sum_{i=1}^n c_i \mathbf{x}_i = \sum_{i=1}^n c_i A \mathbf{x}_i = \sum_{i=1}^n c_i \lambda_i \mathbf{x}_i$$

$$\mathbf{0} = \lambda_n \sum_{i=1}^n c_i \mathbf{x}_i = \sum_{i=1}^n c_i \lambda_n \mathbf{x}_i$$

Subtract those two equations so that the coefficient of \mathbf{x}_n is $\lambda_n - \lambda_n$ and \mathbf{x}_n is gone. Now multiply by A and λ_{n-1} and subtract. \mathbf{x}_{n-1} is gone.

$$\text{We reach } (\lambda_1 - \lambda_2) \dots (\lambda_1 - \lambda_n) c_1 \mathbf{x}_1 = \mathbf{0} \implies c_1 = 0$$

Similarly, $c_i = 0$ for all i 's and the eigenvectors are independent. ■

Corollary 1: Diagonalizable Matrix

An n by n matrix with n distinct eigenvalues must be diagonalizable.

Symmetric Matrices

Symmetric matrices ($A^T = A$) are very important. And they have some very useful properties. And their ability to diagonalize is extraordinary. At first glance, their symmetry reveals:

$$X \Lambda X^{-1} = A = A^T = (X \Lambda X^{-1})^T = (X^{-1})^T \Lambda X^T$$

Does this suggest that $X^T = X^{-1}$? If so, then $X^T X = I$ and each eigenvector in X is orthogonal to the other eigenvectors. Even better, we could write $A = Q \Lambda Q^T$. It is much more efficient to compute a transpose than an inverse; so this is a big deal by itself. But we will see much greater benefits as well.

Lemma 3: Properties Of Eigenvalues And Eigenvectors For Symmetric Matrix

For a real, symmetric matrix:

- ① The eigenvalues are real
- ② The eigenvectors corresponding to distinct eigenvalues are orthogonal

Proof:

For ①, suppose there exists a complex λ such that $S\mathbf{x} = \lambda\mathbf{x}$. Then

$$S\mathbf{x} \stackrel{\phi}{=} \lambda\mathbf{x} \implies S\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}} \implies \mathbf{x}^H S = (S\bar{\mathbf{x}})^T \stackrel{\psi}{=} (\bar{\lambda}\bar{\mathbf{x}})^T = \mathbf{x}^H \bar{\lambda}$$

Now multiply ϕ on the left by \mathbf{x}^H and multiply ψ on the right by \mathbf{x} :

$$\mathbf{x}^H S \mathbf{x} = \mathbf{x}^H \lambda \mathbf{x} \quad \mathbf{x}^H S \mathbf{x} = \mathbf{x}^H \bar{\lambda} \mathbf{x}$$

So $\mathbf{x}^H \lambda \mathbf{x} = \mathbf{x}^H \bar{\lambda} \mathbf{x}$. Hence $\lambda = \bar{\lambda}$ and λ is real.

For ②, let $S\mathbf{x} = \lambda_1\mathbf{x}$ and $S\mathbf{y} = \lambda_2\mathbf{y}$ where $\lambda_1 \neq \lambda_2$.

$$\mathbf{x}^T \lambda_1 \mathbf{y} = (\lambda_1 \mathbf{x})^T \mathbf{y} = (S\mathbf{x})^T \mathbf{y} = \mathbf{x}^T S^T \mathbf{y} = \mathbf{x}^T S \mathbf{y} = \mathbf{x}^T \lambda_2 \mathbf{y}$$

So $\mathbf{x}^T \mathbf{y} = 0$. ■

Corollary 2: Limited Spectral Theorem

Given a real, symmetric matrix S with distinct eigenvalues:

- ① Its eigenvector matrix Q is orthogonal and invertible and $Q^{-1} = Q^T$.
- ② It's diagonalizable: $S = Q\Lambda Q^{-1} = Q\Lambda Q^T = \sum_{i=1}^n \lambda_i \mathbf{q}_i \mathbf{q}_i^T$

For ①, since its eigenvectors are orthogonal, they are linearly independent and Q is invertible. We can normalize the eigenvectors so that Q has orthonormal columns and $Q^T = Q^{-1}$.

For ②, we assumed distinct eigenvalues, so Corollary 1 ensures it is diagonalizable with $S = Q\Lambda Q^{-1}$. Using ①, we have $Q^{-1} = Q^T$ so that $S = Q\Lambda Q^{-1} = Q\Lambda Q^T$. ■

You might be thinking... OK, that's all great. $Q^T = Q^{-1}$ makes the diagonalization more efficient to compute. Whoop-dee-doo. Didn't you say that symmetric matrices have some extraordinary ability to diagonalize? We already knew from Corollary 1 that any matrix with distinct eigenvalues is diagonalizable. Is the "distinct eigenvalues" disclaimer superfluous for symmetric matrices? Are you suggesting that any symmetric matrix is diagonalizable?

Spectral Theorem

Every symmetric matrix has the factorization $S = Q\Lambda Q^T$ with real eigenvalues in Λ and orthonormal eigenvectors in the columns of Q .

$$S = Q\Lambda Q^T = Q\Lambda Q^{-1} = \sum_{i=1}^n \lambda_i \mathbf{q}_i \mathbf{q}_i^T \text{ and } Q^{-1} = Q^T$$

Proof:

Many proofs exist. [Inductive Proof](#). [Another Proof](#). It is also possible to prove that you can get an orthonormal set of eigenvectors for any symmetric matrix, even for the case of repeated eigenvalues (basically apply Gram-Schmidt to the eigenspaces of repeated eigenvalues). So that's another proof. Later I will present another proof using the Rayleigh Quotient.

$A^T A$ and AA^T Properties

Let A be a real m by n matrix of rank r . Then $A^T A$ and AA^T

- ① are symmetric
- ② are positive semi-definite (and positive definite if A is square and non-singular)
- ③ share the same nullspace as A and A^T , respectively.
- ④ share rank r .
- ⑤ share the same determinant ($=|A|^2$ if A is square, haha)
- ⑥ share the same trace
- ⑦ share the same eigenvalues
- ⑧ share related eigenvectors: $\mathbf{v} = A^T \mathbf{u}$ and $\mathbf{u} = A \mathbf{v}$
- ⑨ are diagonalizable with orthonormal eigenvectors, $Q^{-1} = Q^T$, and $Q\Lambda Q^T = \sum_{i=1}^n \lambda_i \mathbf{q}_i \mathbf{q}_i^T$
- ⑩ have r positive eigenvalues

Proof:

To show ①:

$$(A^T A)^T = A^T (A^T)^T = A^T A$$

For ②, we already showed symmetry. We must also show non-negative eigenvalues. It is sufficient to show that $\mathbf{x}^T A^T A \mathbf{x} \geq 0$ for any vector \mathbf{x} (the energy-based definition):

$$\mathbf{x}^T A^T A \mathbf{x} = (A \mathbf{x})^T A \mathbf{x} = \|A \mathbf{x}\|^2 \geq 0$$

$$\mathbf{y}^T A A^T \mathbf{y} = (A^T \mathbf{y})^T A^T \mathbf{y} = \|A^T \mathbf{y}\|^2 \geq 0$$

In particular, if A is non-singular, then $A^T A$ and AA^T are positive definite.

For ③, suppose $\mathbf{x} \in N(A)$ so that $A \mathbf{x} = \mathbf{0}$. Then $A^T A \mathbf{x} = A^T \mathbf{0} = \mathbf{0}$. So $\mathbf{x} \in N(A^T A)$. Conversely, suppose that $\mathbf{x} \in N(A^T A)$ so that $A^T A \mathbf{x} = \mathbf{0}$. Left multiplying by \mathbf{x}^T , we get

$$0 = \mathbf{x}^T \mathbf{0} = \mathbf{x}^T A^T A \mathbf{x} = (A \mathbf{x})^T (A \mathbf{x}) = \|A \mathbf{x}\|^2$$

Hence $A \mathbf{x} = \mathbf{0}$ and $\mathbf{x} \in N(A)$.

We have shown that $N(A) \subset N(A^T A)$ and $N(A^T A) \subset N(A)$. Hence $N(A^T A) = N(A)$.

For AA^T , suppose $\mathbf{y} \in N(A^T)$ so that $A^T \mathbf{y} = \mathbf{0}$. Then $AA^T \mathbf{y} = A\mathbf{0} = \mathbf{0}$. So $\mathbf{y} \in N(AA^T)$. Conversely, suppose that $\mathbf{y} \in N(AA^T)$ so that $AA^T \mathbf{y} = \mathbf{0}$. Left multiplying by \mathbf{y}^T , we get

$$0 = \mathbf{y}^T \mathbf{0} = \mathbf{y}^T AA^T \mathbf{y} = (A^T \mathbf{y})^T (A^T \mathbf{y}) = \|A^T \mathbf{y}\|^2$$

Hence $A^T \mathbf{y} = \mathbf{0}$ and $\mathbf{y} \in N(A^T)$.

We have shown that $N(A^T) \subset N(AA^T)$ and $N(AA^T) \subset N(A^T)$. Hence $N(AA^T) = N(A^T)$.

To prove (4), (3) $\implies \dim(N(A^T A)) = \dim(N(A)) = n - r$. Hence $\dim(C(A^T A)) = n - (n - r) = r$ and the rank of $A^T A$ is r . Similarly, we can show that r is the rank of AA^T .

For (5):

$$|A^T A| = |A^T| |A| = |A| |A^T| = |AA^T|$$

And if A is square: $|A^T A| = |AA^T| = |A^T| |A| = |A| |A| = |A|^2$

For (7) and (8), let λ be an eigenvalue of AA^T with eigenvector \mathbf{u} . And let $\mathbf{v} = A^T \mathbf{u}$. Then

$$A\mathbf{v} = AA^T \mathbf{u} = \lambda \mathbf{u}$$

$$\implies A^T A\mathbf{v} = A^T \lambda \mathbf{u} = \lambda A^T \mathbf{u} = \lambda \mathbf{v}$$

Hence λ is an eigenvalue of $A^T A$ with eigenvector $\mathbf{v} = A^T \mathbf{u}$. On the other hand, let λ be an eigenvalue of $A^T A$ with eigenvector \mathbf{v} . And let $\mathbf{u} = A\mathbf{v}$. Then

$$A^T \mathbf{u} = A^T A\mathbf{v} = \lambda \mathbf{v}$$

$$\implies AA^T \mathbf{u} = A\lambda \mathbf{v} = \lambda A\mathbf{v} = \lambda \mathbf{u}$$

Hence λ is an eigenvalue of AA^T with eigenvector $\mathbf{u} = A\mathbf{v}$.

(9) follows directly from the Spectral Theorem.

For (10), non-negativity of eigenvalues follows from positive semi-definiteness proven in (2). But now we want to know the number of positive eigenvalues vs. the number of zero eigenvalues. (9) shows $A^T A$ is diagonalizable. So, for every eigenvalue of $A^T A$, the geometric multiplicity equals the algebraic multiplicity. So n minus the dimension of the eigenspace of zero will equal the number of positive eigenvalues. In (3), we showed that $A^T A$ has the same nullspace as A . But the nullspace of $A^T A$ is just the eigenspace of zero for $A^T A$. So

Number of positive eigenvalues for $A^T A = n - \text{dimension of the eigenspace of zero for } A^T A$

$$= n - \dim(N(A^T A)) = n - \dim(N(A))$$

$$= n - (n - r) = r$$

A similar argument works for AA^T . ■

Theorem: Singular Value Decomposition

Let A be an m by n matrix of rank r . There exists a diagonalization $A = U\Sigma V^T$ such that:

- ① The first r columns of U form an orthonormal basis for the column space of A .
- ② The first r columns of V form an orthonormal basis for the row space of A .
- ③ The first r elements of the diagonal of Σ are square roots of eigenvalues of $A^T A$ (and equivalently AA^T).
- ④ The remaining $m - r$ columns of U form an orthonormal basis for the left nullspace of A .
- ⑤ The remaining $n - r$ columns of V form an orthonormal basis for the nullspace of A .
- ⑥ $A = U\Sigma V^T$ separates A into r rank 1 matrices $\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$.
- ⑦ $A^T A = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$ and the right singular vectors of A are eigenvectors of $A^T A$ with eigenvalues $\sigma_i^2 = \lambda_i$.
- ⑧ $AA^T = \sum_{i=1}^r \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^T$ and the left singular vectors of A are eigenvectors of AA^T with eigenvalues $\sigma_i^2 = \lambda_i$.
- ⑨ The columns of U form an orthonormal basis for \mathbb{R}^m .
- ⑩ The columns of V form an orthonormal basis for \mathbb{R}^n .

There are at least a couple of ways to prove the SVD Theorem. I present two proofs that I find to be very instructive. "All at once" and "one at a time".

SVD Proof: All At Once

Suppose $r \leq n \leq m$. The proof for $m < n$ is similar.

We choose the \mathbf{v}_i 's (AKA the right singular vectors) to be orthonormal eigenvectors of $A^T A$ with corresponding eigenvalues λ_i 's. **$A^T A$ and AA^T Properties (9)** guarantees we can get n such pairs of eigenvalues and orthonormal eigenvectors. We know from **$A^T A$ and AA^T Properties (10)** that there are r positive eigenvalues and $n - r$ zero eigenvalues. Index the eigenvalues and eigenvectors so that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ and $\lambda_{r+1} = \dots = \lambda_n = 0$. Define the singular values σ_i as

$$\sigma_i \equiv \|A\mathbf{v}_i\| = \sqrt{(A\mathbf{v}_i)^T A\mathbf{v}_i} = \sqrt{\mathbf{v}_i^T A^T A \mathbf{v}_i} = \sqrt{\mathbf{v}_i^T \lambda_i \mathbf{v}_i} = \sqrt{\lambda_i} \sqrt{\mathbf{v}_i^T \mathbf{v}_i} = \sqrt{\lambda_i} \|\mathbf{v}_i\| = \sqrt{\lambda_i}$$

and, for $\sigma_i > 0$, define $\mathbf{u}_i \equiv \frac{A\mathbf{v}_i}{\sigma_i}$. Then

$$\mathbf{u}_i^T \mathbf{u}_j = \left(\frac{A\mathbf{v}_i}{\sigma_i} \right)^T \left(\frac{A\mathbf{v}_j}{\sigma_j} \right) = \frac{\mathbf{v}_i^T A^T A \mathbf{v}_j}{\sigma_i \sigma_j} = \frac{\mathbf{v}_i^T \lambda_j \mathbf{v}_j}{\sigma_i \sigma_j} = 0 \quad \text{for } i = 1, 2, \dots, r$$

The \mathbf{u}_i 's inherit orthogonality from the \mathbf{v}_i 's. This is vital to SVD. And the \mathbf{u}_i 's are also unit vectors:

$$\|\mathbf{u}_i\| = \left\| \frac{A\mathbf{v}_i}{\sigma_i} \right\| = \frac{\|A\mathbf{v}_i\|}{\sigma_i} = 1.$$

Now we wish to extend the \mathbf{u}_i 's (AKA the left singular vectors) to m vectors. And we wish to extend the singular values to a count of m . Choose $\sigma_{n+1} = \dots = \sigma_m = 0$. Any choice for the remaining $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ vectors will trivially satisfy $A\mathbf{v}_i = \sigma_i \mathbf{u}_i$ since $\sigma_{r+1} = \dots = \sigma_m = 0$. We choose some orthonormal basis for the nullspace of A^T , which has dimension $m - r$. So choosing a basis from the nullspace of A^T will give us $m - r$ vectors, which is the number of vectors we want for $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$.

Define $U_{m \times m} \equiv [\mathbf{u}_1 \ \dots \ \mathbf{u}_m]$, $V_{n \times n} \equiv [\mathbf{v}_1 \ \dots \ \mathbf{v}_n]$, and

$$\Sigma_{m \times n} \equiv \begin{bmatrix} \sigma_1 & & & & & \\ & \sigma_2 & & & & \\ & & \ddots & & & \\ & & & \sigma_r & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix}$$

We now have all the pieces in place. Let's start proving stuff.

For ①, we already showed that the first r left singular vectors $(\mathbf{u}_1, \dots, \mathbf{u}_r)$ are orthonormal. So they form an orthonormal basis for some r dimensional subspace of \mathbb{R}^m . But is that subspace the column space of A ? Well, we chose $\mathbf{u}_i = \frac{A\mathbf{v}_i}{\sigma_i} = A\mathbf{w}_i$ where $\mathbf{w}_i = \frac{1}{\sigma_i}\mathbf{v}_i$. So, by definition, the $\mathbf{u}_1, \dots, \mathbf{u}_r$ are in the column space of A .

For ②, we chose orthonormal \mathbf{v}_i 's. And for $i \leq r$, we have

$$\lambda_i \mathbf{v}_i = A^T A \mathbf{v}_i = A^T \sigma_i \mathbf{u}_i \iff \mathbf{v}_i = A^T \mathbf{w}_i \quad \text{where } \mathbf{w}_i = \frac{\sigma_i}{\lambda_i} \mathbf{u}_i$$

Hence the $\mathbf{v}_1, \dots, \mathbf{v}_r$'s are in the row space of A and form an orthonormal basis for the row space of A .

For ③, you can clearly see this.

For ④, note that we chose $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ as an orthonormal basis for the nullspace of A^T .

For ⑤, we chose orthonormal \mathbf{v}_i 's. And we indexed the \mathbf{v}_i 's and λ_i 's such that $\lambda_{r+1} = \dots = \lambda_n = 0$. Hence $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$ are in the zero eigenspace of $A^T A$, which is the nullspace of $A^T A$, which is the nullspace of A , by **$A^T A$ and AA^T Properties (3)**.

For ⑥, first let's check that $A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$:

$$U \Sigma V^T = \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & & \\ & \sigma_2 & & & & \\ & & \ddots & & & \\ & & & \sigma_r & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix}$$

$$\begin{aligned}
&= \begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_m \end{bmatrix} \begin{bmatrix} \sigma_1 v_{1,1} & \sigma_1 v_{1,2} & \dots & \sigma_1 v_{1,r} & \sigma_1 v_{1,r+1} & \dots & \sigma_1 v_{1,n} \\ \sigma_2 v_{2,1} & \sigma_2 v_{2,2} & \dots & \sigma_2 v_{2,r} & \sigma_2 v_{2,r+1} & \dots & \sigma_2 v_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \dots & \vdots \\ \sigma_r v_{r,1} & \sigma_r v_{r,2} & \dots & \sigma_r v_{r,r} & \sigma_r v_{r,r+1} & \dots & \sigma_r v_{r,n} \\ 0 & & & & 0 & & \\ & \ddots & & & & \ddots & \\ & & 0 & 0 & & & 0 \end{bmatrix} \\
&= \begin{bmatrix} \sum_{i=1}^r \sigma_i u_{i,1} v_{i,1} & \sum_{i=1}^r \sigma_i u_{i,1} v_{i,2} & \dots & \sum_{i=1}^r \sigma_i u_{i,1} v_{i,r} & \sum_{i=1}^r \sigma_i u_{i,1} v_{i,r+1} & \dots & \sum_{i=1}^r \sigma_i u_{i,1} v_{i,n} \\ \sum_{i=1}^r \sigma_i u_{i,2} v_{i,1} & \sum_{i=1}^r \sigma_i u_{i,2} v_{i,2} & \dots & \sum_{i=1}^r \sigma_i u_{i,2} v_{i,r} & \sum_{i=1}^r \sigma_i u_{i,2} v_{i,r+1} & \dots & \sum_{i=1}^r \sigma_i u_{i,2} v_{i,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \dots & \vdots \\ \sum_{i=1}^r \sigma_i u_{i,r} v_{i,1} & \sum_{i=1}^r \sigma_i u_{i,r} v_{i,2} & \dots & \sum_{i=1}^r \sigma_i u_{i,r} v_{i,r} & \sum_{i=1}^r \sigma_i u_{i,r} v_{i,r+1} & \dots & \sum_{i=1}^r \sigma_i u_{i,r} v_{i,n} \\ \sum_{i=1}^r \sigma_i u_{i,r+1} v_{i,1} & \sum_{i=1}^r \sigma_i u_{i,r+1} v_{i,2} & \dots & \sum_{i=1}^r \sigma_i u_{i,r+1} v_{i,r} & \sum_{i=1}^r \sigma_i u_{i,r+1} v_{i,r+1} & \dots & \sum_{i=1}^r \sigma_i u_{i,r+1} v_{i,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^r \sigma_i u_{i,m} v_{i,1} & \sum_{i=1}^r \sigma_i u_{i,m} v_{i,2} & \dots & \sum_{i=1}^r \sigma_i u_{i,m} v_{i,r} & \sum_{i=1}^r \sigma_i u_{i,m} v_{i,r+1} & \dots & \sum_{i=1}^r \sigma_i u_{i,m} v_{i,n} \end{bmatrix} = \Phi
\end{aligned}$$

To show that $\Phi = A$, we must show that $A[j, k] = \sum_{i=1}^r \sigma_i u_{i,j} v_{i,k}$. Recall that $\mathbf{u}_i = \frac{\mathbf{A} \mathbf{v}_i}{\sigma_i}$. Hence

$$u_{i,j} = \frac{\mathbf{A}_j \mathbf{v}_i}{\sigma_i} \text{ and}$$

$$\sum_{i=1}^r \sigma_i u_{i,j} v_{i,k} = \sum_{i=1}^r \sigma_i \left(\frac{\mathbf{A}_j \mathbf{v}_i}{\sigma_i} \right) v_{i,k} = \mathbf{A}_j \sum_{i=1}^r \mathbf{v}_i v_{i,k}$$

So we must show that $\sum_{i=1}^r \mathbf{v}_i v_{i,k}$ is the k^{th} standard unit vector in \mathbb{R}^n :

$$\begin{aligned}
\sum_{i=1}^r \mathbf{v}_i v_{i,k} &= \sum_{i=1}^r \begin{bmatrix} v_{i,1} \\ v_{i,2} \\ \vdots \\ v_{i,k} \\ \vdots \\ v_{i,n} \end{bmatrix} v_{i,k} = \sum_{i=1}^r \begin{bmatrix} v_{i,1} v_{i,k} \\ v_{i,2} v_{i,k} \\ \vdots \\ v_{i,k}^2 \\ \vdots \\ v_{i,n} v_{i,k} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^r v_{i,1} v_{i,k} \\ \sum_{i=1}^r v_{i,2} v_{i,k} \\ \vdots \\ \sum_{i=1}^r v_{i,k}^2 \\ \vdots \\ \sum_{i=1}^r v_{i,n} v_{i,k} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}
\end{aligned}$$

The last equation follows from the fact that the transpose of an orthogonal matrix (V) is also an orthogonal matrix. Hence $\Phi = A$.

Now let's look at the dyadic sum $\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. It is clear that each term in that sum is rank 1. It is instructive to write out the first couple of terms in this sum:

$$\begin{aligned}
\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T &= \sigma_1 \begin{bmatrix} u_{1,1} \\ u_{1,2} \\ \vdots \\ u_{1,m} \end{bmatrix} \begin{bmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \end{bmatrix} = \begin{bmatrix} \sigma_1 u_{1,1} v_{1,1} & \sigma_1 u_{1,1} v_{1,2} & \dots & \sigma_1 u_{1,1} v_{1,n} \\ \sigma_1 u_{1,2} v_{1,1} & \sigma_1 u_{1,2} v_{1,2} & \dots & \sigma_1 u_{1,2} v_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_1 u_{1,m} v_{1,1} & \sigma_1 u_{1,m} v_{1,2} & \dots & \sigma_1 u_{1,m} v_{1,n} \end{bmatrix} \\
\sigma_2 \mathbf{u}_2 \mathbf{v}_2^T &= \sigma_2 \begin{bmatrix} u_{2,1} \\ u_{2,2} \\ \vdots \\ u_{2,m} \end{bmatrix} \begin{bmatrix} v_{2,1} & v_{2,2} & \dots & v_{2,n} \end{bmatrix} = \begin{bmatrix} \sigma_2 u_{2,1} v_{2,1} & \sigma_2 u_{2,1} v_{2,2} & \dots & \sigma_2 u_{2,1} v_{2,n} \\ \sigma_2 u_{2,2} v_{2,1} & \sigma_2 u_{2,2} v_{2,2} & \dots & \sigma_2 u_{2,2} v_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_2 u_{2,m} v_{2,1} & \sigma_2 u_{2,m} v_{2,2} & \dots & \sigma_2 u_{2,m} v_{2,n} \end{bmatrix}
\end{aligned}$$

And summing them:

$$\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T = \begin{bmatrix} \sum_{i=1}^2 \sigma_i u_{i,1} v_{i,1} & \sum_{i=1}^2 \sigma_i u_{i,1} v_{i,2} & \cdots & \sum_{i=1}^2 \sigma_i u_{i,1} v_{i,n} \\ \sum_{i=1}^2 \sigma_i u_{i,2} v_{i,1} & \sum_{i=1}^2 \sigma_i u_{i,2} v_{i,2} & \cdots & \sum_{i=1}^2 \sigma_i u_{i,2} v_{i,n} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^2 \sigma_i u_{i,m} v_{i,1} & \sum_{i=1}^2 \sigma_i u_{i,m} v_{i,2} & \cdots & \sum_{i=1}^2 \sigma_i u_{i,m} v_{i,n} \end{bmatrix}$$

And it is clear that $\Phi = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$.

For (7):

$$A^T A = \left(\sum_{i=1}^r \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right) = \sum_{i,j=1}^r \sigma_i \sigma_j \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{v}_j^T = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T$$

because the orthonormality of the \mathbf{u}_i 's means $\mathbf{u}_i^T \mathbf{u}_j = 0$ for $i \neq j$ and $\mathbf{u}_i^T \mathbf{u}_j = 1$ for $i = j$.

$$A^T A \mathbf{v}_j = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \right) \mathbf{v}_j = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \mathbf{v}_j = \sigma_j^2 \mathbf{v}_j$$

thanks again to orthonormality.

For (8):

$$A A^T = \left(\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^r \sigma_j \mathbf{v}_j \mathbf{u}_j^T \right) = \sum_{i,j=1}^r \sigma_i \sigma_j \mathbf{u}_i \mathbf{v}_i^T \mathbf{v}_j \mathbf{u}_j^T = \sum_{i=1}^r \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^T$$

$$A A^T \mathbf{u}_j = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^T \right) \mathbf{u}_j = \sum_{i=1}^r \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^T \mathbf{u}_j = \sigma_j^2 \mathbf{u}_j$$

For (9), there are m orthonormal $m \times 1$ column vectors in U . Or are there? We never showed that each of the $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ is orthogonal to each of the $\mathbf{u}_1, \dots, \mathbf{u}_r$. But (1) showed that the $\mathbf{u}_1, \dots, \mathbf{u}_r$ are in the column space of A . And (4) showed that the $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$ are in the nullspace of A^T . And these two spaces are orthogonal complements. Hence the $\mathbf{u}_1, \dots, \mathbf{u}_m$ are orthonormal and they form an orthonormal basis for \mathbb{R}^m .

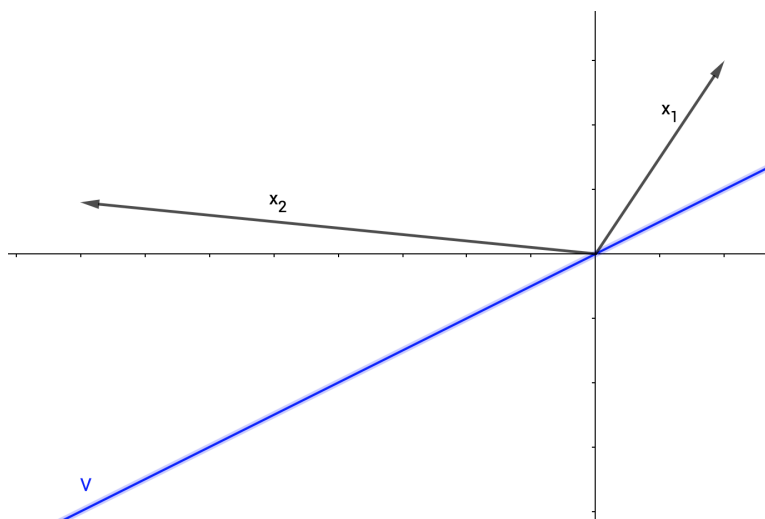
For (10), there are n orthonormal $n \times 1$ columns in V . This forms an orthonormal basis for \mathbb{R}^n . ■

You can see why I called this proof the "all at once" proof. We got all the singular values and all the right singular vectors in one shot.

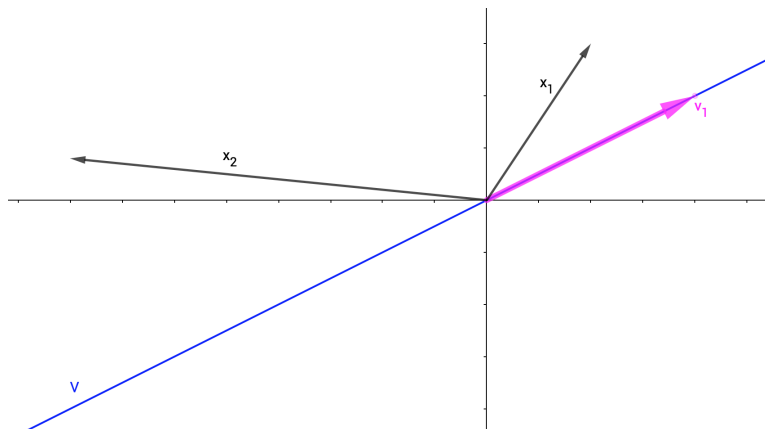
There is a different way to look at SVD. This will lead to a restatement of SVD. It will also lead to the "one at a time" proof of SVD that I mentioned before. This approach, sometimes called the **best least squares fit**, starts with a question: Say we have a set $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ of m vectors in \mathbb{R}^n . Can we find the best k -dimensional subspace V of \mathbb{R}^n with respect to this set of points? We define "best" to mean the shortest distance from X to the subspace V . More specifically, we define "best" to mean the minimum sum of the squares of the perpendicular distances from the points in X to the subspace V .

I know of at least two great resources that treat SVD as a **best least squares fit** problem. One resource is a pair of blog posts by Jeremy Kun ([SVD 1](#) and [SVD 2](#)). I highly recommend both posts. He provides great intuition into a number of topics in these two posts. The other resource is [here](#). But I'm sure there are other excellent treatments of SVD as a *best least squares fit* problem out there.

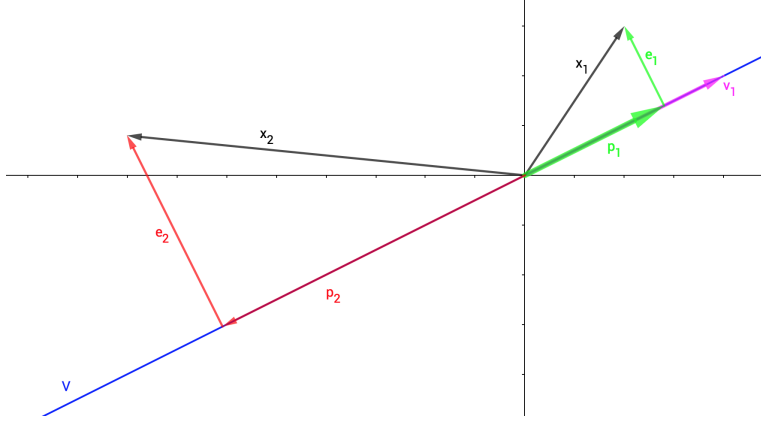
OK, let's start with some simple pictures. Once we get the geometry, the best least squares fit approach to SVD and the "one at a time" proof follow very quickly. Let's say that $m = n = 2$ and let's say $k = 1$. So $X = \{x_1, x_2\}$ consists of two vectors in \mathbb{R}^2 . For the moment, the pictured V can be any line through the origin in \mathbb{R}^2 . The pictured V is not yet the "best-fit" line:



If we want to minimize the sum of the squares of the perpendicular distances from the points in X to the subspace V , then we will need a basis for V . Let's call $\{v_1\}$ an orthonormal basis for (the yet non-optimal subspace) V .



And now we want to see the projections of X onto V :



\mathbf{p}_1 is the projection of \mathbf{x}_1 onto \mathbf{v}_1 and \mathbf{p}_2 is the projection of \mathbf{x}_2 onto \mathbf{v}_1 . $\mathbf{e}_1 = \mathbf{x}_1 - \mathbf{p}_1$ is the error from the projection of \mathbf{x}_1 onto \mathbf{v}_1 and $\mathbf{e}_2 = \mathbf{x}_2 - \mathbf{p}_2$ is the error from the projection of \mathbf{x}_2 onto \mathbf{v}_1 . Recall that, for a given \mathbf{v}_1 , the error \mathbf{e}_i is minimized when we put $\mathbf{p}_i = \frac{\mathbf{v}_1^T \mathbf{x}_i}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1$. It is easy to see and compute that \mathbf{e}_1 and \mathbf{e}_2 are perpendicular to \mathbf{v}_1 :

$$\begin{aligned} \mathbf{e}_i \cdot \mathbf{v}_1 &= \mathbf{x}_i^T \mathbf{v}_1 - \mathbf{p}_i^T \mathbf{v}_1 = \mathbf{x}_i^T \mathbf{v}_1 - \left(\frac{\mathbf{v}_1^T \mathbf{x}_i}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1 \right)^T \mathbf{v}_1 \\ &= \mathbf{x}_i^T \mathbf{v}_1 - \left(\frac{\mathbf{x}_i^T \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1^T \right) \mathbf{v}_1 = \mathbf{x}_i^T \mathbf{v}_1 - \mathbf{x}_i^T \mathbf{v}_1 = 0 \end{aligned}$$

Similarly, \mathbf{e}_i is perpendicular to \mathbf{p}_i :

$$\begin{aligned} \mathbf{e}_i \cdot \mathbf{p}_i &= \mathbf{x}_i^T \mathbf{p}_i - \mathbf{p}_i^T \mathbf{p}_i = \mathbf{x}_i^T \mathbf{p}_i - \left(\frac{\mathbf{v}_1^T \mathbf{x}_i}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1 \right)^T \mathbf{p}_i \\ &= \mathbf{x}_i^T \mathbf{p}_i - \left(\frac{\mathbf{x}_i^T \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1^T \right) \mathbf{p}_i = \mathbf{x}_i^T \mathbf{p}_i - \frac{\mathbf{x}_i^T \mathbf{v}_1}{\|\mathbf{v}_1\|^2} \mathbf{p}_i^T \mathbf{v}_1 \\ &= \mathbf{x}_i^T \mathbf{p}_i - \mathbf{x}_i^T \mathbf{v}_1 \mathbf{p}_i^T \mathbf{v}_1 = \mathbf{x}_i^T \mathbf{p}_i - \mathbf{x}_i^T \mathbf{p}_i \mathbf{v}_1^T \mathbf{v}_1 = \mathbf{x}_i^T \mathbf{p}_i - \mathbf{x}_i^T \mathbf{p}_i = 0 \end{aligned}$$

Notice that

$$0 = \mathbf{e}_i \cdot \mathbf{p}_i = \mathbf{x}_i^T \mathbf{p}_i - \mathbf{p}_i^T \mathbf{p}_i = \mathbf{p}_i^T \mathbf{x}_i - \mathbf{p}_i^T \mathbf{p}_i$$

So that

$$\mathbf{x}_i^T \mathbf{p}_i = \mathbf{p}_i^T \mathbf{p}_i \text{ and } \mathbf{p}_i^T \mathbf{x}_i = \mathbf{p}_i^T \mathbf{p}_i$$

Hence, for a given \mathbf{v}_1 , we have:

$$\begin{aligned} \|\mathbf{e}_i\|^2 + \|\mathbf{p}_i\|^2 &= \mathbf{e}_i^T \mathbf{e}_i + \mathbf{p}_i^T \mathbf{p}_i = (\mathbf{x}_i - \mathbf{p}_i)^T (\mathbf{x}_i - \mathbf{p}_i) + \mathbf{p}_i^T \mathbf{p}_i \\ &= (\mathbf{x}_i^T - \mathbf{p}_i^T) (\mathbf{x}_i - \mathbf{p}_i) + \mathbf{p}_i^T \mathbf{p}_i = \mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{p}_i - \mathbf{p}_i^T \mathbf{x}_i + \mathbf{p}_i^T \mathbf{p}_i + \mathbf{p}_i^T \mathbf{p}_i = \|\mathbf{x}_i\|^2 \end{aligned}$$

Succintly, for a given \mathbf{v}_1 , we have the Pythagorean Theorem:

$$\|\mathbf{e}_i\|^2 + \|\mathbf{p}_i\|^2 = \|\mathbf{x}_i\|^2$$

Now remember that we haven't yet selected the "best-fit" subspace V nor its basis $\{\mathbf{v}_1\}$. The objective is to find the "optimal" 1-dimensional V such that we minimize the sum of the squares of the perpendicular distances from the points in X to any 1-dimensional subspace. In our simple pictures, the sum of the squares of the perpendicular distances is $\|\mathbf{e}_1\|^2 + \|\mathbf{e}_2\|^2$:

$$\|\mathbf{e}_1\|^2 + \|\mathbf{e}_2\|^2 = \|\mathbf{x}_1\|^2 + \|\mathbf{x}_2\|^2 - \|\mathbf{p}_1\|^2 - \|\mathbf{p}_2\|^2$$

So how can we choose V and its basis $\{\mathbf{v}_1\}$ so as to minimize $\|\mathbf{e}_1\|^2 + \|\mathbf{e}_2\|^2$? Well \mathbf{x}_1 and \mathbf{x}_2 are given and thus fixed. But $\mathbf{p}_1 \equiv \frac{\mathbf{v}_1^T \mathbf{x}_1}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1$ and $\mathbf{p}_2 \equiv \frac{\mathbf{v}_1^T \mathbf{x}_2}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1$ are dependent on \mathbf{v}_1 , when we select \mathbf{p}_1 and \mathbf{p}_2 "correctly". Hence we can minimize the sum of the squares of the perpendicular distances by choosing \mathbf{v}_1 to maximize $\|\mathbf{p}_1\|^2 + \|\mathbf{p}_2\|^2$. More generally, we can minimize the sum of the squares of the perpendicular distances between X and V by maximizing the sum of the squared lengths of the projections of X onto V . That is, we have converted our original minimization problem into a maximization problem:

$$\begin{aligned} \mathbf{v}_1 &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} (\|\mathbf{p}_1\|^2 + \|\mathbf{p}_2\|^2) = \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \left(\sum_{i=1}^2 \|\mathbf{p}_i\|^2 \right) \\ &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \left(\sum_{i=1}^2 \left\| \frac{\mathbf{v}^T \mathbf{x}_i}{\mathbf{v}^T \mathbf{v}} \mathbf{v} \right\|^2 \right) = \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \left(\sum_{i=1}^2 \left\| \frac{\mathbf{v}^T \mathbf{x}_i}{1} \mathbf{v} \right\|^2 \right) \\ &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \left(\sum_{i=1}^2 \left\| (\mathbf{v}^T \mathbf{x}_i) \mathbf{v} \right\|^2 \right) = \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \left(\sum_{i=1}^2 (\mathbf{v}^T \mathbf{x}_i)^2 \|\mathbf{v}\|^2 \right) \\ &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \sum_{i=1}^2 (\mathbf{v} \cdot \mathbf{x}_i)^2 \end{aligned}$$

But we can write this more compactly by putting the \mathbf{x}_i 's into the rows of a matrix, call it A :

$$A = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \end{bmatrix}$$

So that

$$A\mathbf{v} = \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 \cdot \mathbf{v} \\ \mathbf{x}_2 \cdot \mathbf{v} \end{bmatrix}$$

And

$$\|A\mathbf{v}\|^2 = (\mathbf{x}_1 \cdot \mathbf{v})^2 + (\mathbf{x}_2 \cdot \mathbf{v})^2$$

And

$$\begin{aligned} \mathbf{v}_1 &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \sum_{i=1}^2 (\mathbf{x}_i \cdot \mathbf{v})^2 = \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\|^2 \\ &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\| \end{aligned}$$

The last equation follows from the fact the squaring function is a monotonically increasing function.

Now let's take this up a notch. Suppose $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ with $\mathbf{x}_i \in \mathbb{R}^n$. We want to find the plane V in \mathbb{R}^n that minimizes the sum of the squares of the perpendicular distances. That is, the objective is to find an orthonormal basis $\{\mathbf{v}_1, \mathbf{v}_2\}$ with $\mathbf{v}_i \in \mathbb{R}^n$ such that we minimize the sum of the squares of the perpendicular distances between $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ and $\{\mathbf{v}_1, \mathbf{v}_2\}$.

Or even more generally, let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ with $\mathbf{x}_i \in \mathbb{R}^n$. The objective is to find some k -dimensional subspace of \mathbb{R}^n , call it V , with orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ such that we minimize the sum of the squares of the perpendicular distances between $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$. Equivalently, the objective is to find some k -dimensional subspace $V \subset \mathbb{R}^n$ with orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ such that we maximize the sum of the squared lengths of the projections of $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ onto $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$.

Well just how are we supposed to do that? Maybe "one at a time"? But we must remember to select orthogonal \mathbf{v}_i 's:

$$\begin{aligned}\mathbf{v}_1 &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\| \\ \mathbf{v}_2 &= \arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|=1} \|A\mathbf{v}\|\end{aligned}$$

For \mathbf{v}_2 , you see that I added a second condition to the arg max function: $\mathbf{v} \perp \mathbf{v}_1$. So this will get us orthonormal vectors. But will it minimize the sum of the squares of the perpendicular distances between $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ and $\{\mathbf{v}_1, \mathbf{v}_2\}$?

Before we answer that most important question, let's officially recognize a few things that you have probably noticed. By taking this "one at a time" approach, we are selecting the right singular vectors, the \mathbf{v}_i 's. We are also selecting the singular values and left singular vectors:

$$\begin{aligned}\sigma_1 &\equiv \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\| \\ \sigma_i &\equiv \max_{\mathbf{v} \perp \mathbf{v}_1, \dots, \mathbf{v} \perp \mathbf{v}_{i-1}, \|\mathbf{v}\|=1} \|A\mathbf{v}\| \\ \mathbf{u}_i &\equiv \frac{A\mathbf{v}_i}{\sigma_i}\end{aligned}$$

It is clear that the singular values are monotonically decreasing: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{k-1} \geq \sigma_k$. This follows directly from our definition for each σ_i . Notice that we specify $\|A\mathbf{v}\|$ (without the square) for the definition of σ_1 and σ_i , rather than $\|A\mathbf{v}\|^2$. For the definition of \mathbf{v}_1 and \mathbf{v}_i , we can use $\|A\mathbf{v}\|$ and $\|A\mathbf{v}\|^2$ interchangeably. The monotonicity of the squaring function guarantees that the same \mathbf{v}_i will satisfy both. But for the definition of the singular value σ_i , we specifically want $\|A\mathbf{v}\|$, without the square.

With these definitions, we can actually prove that our "one at a time" approach is equivalent to the "all at once" approach! Let's do that. If you wish to skip this and continue with the "one at a time approach", [click here](#). First we define $r(\mathbf{x}) = \frac{\mathbf{x}^T S \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ for a symmetric matrix S . This is called the Rayleigh Quotient and it comes in quite handy:

Lemma 4: Eigenvectors & The Rayleigh Quotient

Suppose S is a symmetric matrix and let $r(\mathbf{x}) = \frac{\mathbf{x}^T S \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$. Then $\mathbf{x} \neq \mathbf{0}$ solves $\mathbf{0} = \nabla r(\mathbf{x})$ if and only if \mathbf{x} is an eigenvector of S with eigenvalue $r(\mathbf{x})$.

Proof

First we compute the partials:

$$\begin{aligned}\frac{\partial r(\mathbf{x})}{\partial x_j} &= \frac{\frac{\partial}{\partial x_j}(\mathbf{x}^T S \mathbf{x})}{\mathbf{x}^T \mathbf{x}} - \frac{(\mathbf{x}^T S \mathbf{x}) \frac{\partial}{\partial x_j}(\mathbf{x}^T \mathbf{x})}{(\mathbf{x}^T \mathbf{x})^2} \\ &= \frac{2(S\mathbf{x})_j}{\mathbf{x}^T \mathbf{x}} - \frac{(\mathbf{x}^T S \mathbf{x}) 2x_j}{(\mathbf{x}^T \mathbf{x})^2} \\ &= \frac{2}{\mathbf{x}^T \mathbf{x}} (S\mathbf{x} - r(\mathbf{x})\mathbf{x})_j\end{aligned}$$

So our gradient is $\nabla r(\mathbf{x}) = \frac{2}{\mathbf{x}^T \mathbf{x}} (S\mathbf{x} - r(\mathbf{x})\mathbf{x})$ and we find that

$$\begin{aligned}\mathbf{0} &= \nabla r(\mathbf{x}) = \frac{2}{\mathbf{x}^T \mathbf{x}} (S\mathbf{x} - r(\mathbf{x})\mathbf{x}) \\ \iff \mathbf{0} &= S\mathbf{x} - r(\mathbf{x})\mathbf{x} \\ \iff S\mathbf{x} &= r(\mathbf{x})\mathbf{x}\end{aligned}$$

■

This lemma tells us is that the stationary points of $r(\mathbf{x}) = \frac{\mathbf{x}^T S \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ are eigenvectors for S with eigenvalues equal to $r(\mathbf{x})$ at those points. Hence we can set $S = A^T A$ and get the largest eigenvalue by computing $\lambda_1 = \max_{\mathbf{x}} r(\mathbf{x}) = \max_{\mathbf{x}} \frac{\mathbf{x}^T A^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \max_{\mathbf{x}} \frac{\|A\mathbf{x}\|^2}{\|\mathbf{x}\|^2}$ and

$$\begin{aligned}\lambda_1 &= \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\|^2 = \sigma_1^2 \\ \mathbf{v}_1 &= \arg \max_{\mathbf{v}: \|\mathbf{v}\|=1} \|A\mathbf{v}\|^2\end{aligned}$$

So $\lambda_1 = \sigma_1^2$ is the largest eigenvalue of $A^T A$ with corresponding eigenvector \mathbf{v}_1 . Of course these are the same values we computed above to maximize the sum of the squared lengths of the projections of X onto V . Similarly, lemma 4 says that

$$\begin{aligned}\lambda_i &= \max_{\mathbf{v} \in \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_{i-1})^\perp, \|\mathbf{v}\|=1} \|A\mathbf{v}\|^2 = \sigma_i^2 \\ \mathbf{v}_i &= \arg \max_{\mathbf{v} \in \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_{i-1})^\perp, \|\mathbf{v}\|=1} \|A\mathbf{v}\|^2\end{aligned}$$

is the i^{th} largest eigenpair for $A^T A$ since \mathbf{v}_i is a stationary point of $r(\mathbf{x})$.

So how does this show that the "all at once" approach is equivalent to this "one at a time" approach? Well, in the first few lines of that proof, we picked the right singular vectors to be orthonormal eigenvectors of $A^T A$. And we defined the singular values to be the square roots of the corresponding eigenvalues. Everything else in the proof follows from that. But now we have shown that the \mathbf{v}_i 's and σ_i^2 's defined to maximize the sum of the squared lengths of the projections of X onto V are in fact orthonormal eigenpairs for $A^T A$. Hence equivalency.

Incidentally, this leads directly to another proof of the Spectral Theorem. By iteratively computing the λ_i 's and orthonormal \mathbf{v}_i 's for any given symmetric matrix S using the max, argmax of the Rayleigh Quotient, we are also computing $S = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T = V \Lambda V^T$.

OK, back to our "one at a time" approach. Now you might be asking, can we do this iterative step n times? Well, as long as $\|A\mathbf{v}\| > 0$, then you will have n positive singular values: $\sigma_i > 0$ for $i = 1, \dots, n$. And you will have an orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ for \mathbb{R}^n . As we proved in the "all at once" proof, A has rank n .

But what happens when $\sigma_k > 0$ but $\sigma_{k+1} = 0$ for some $k < n$? This must mean that the rank of A is k . Let's see why that must be the case. Here's what we have: For any $j > k$, any \mathbf{v}_j satisfying $\mathbf{v}_j \perp \mathbf{v}_i$ for all $i = 1, \dots, k$ must also satisfy $\|A\mathbf{v}_j\| = \sigma_j = 0$. Such \mathbf{v}_j 's must also satisfy $A\mathbf{v}_j = \mathbf{0}$. That is, the only remaining \mathbf{v}_j 's that are orthogonal to all $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ must be in the nullspace of A . Since the dimension of \mathbb{R}^n is n , surely we can select $n - k$ such \mathbf{v}_j 's. Said another way, we can extend $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ to an orthonormal basis for \mathbb{R}^n but only with \mathbf{v}_j 's that are in the nullspace of A . Hence the nullspace of A has dimension $n - k$ and hence A has rank k .

In this case where the rank of A is k for $k < n$, you have k positive singular values: $\sigma_i > 0$ for $i = 1, \dots, k$. And you have k right singular vectors $\{\mathbf{v}_i : i = 1, \dots, k\}$. Then you can put $\sigma_i = 0$ for $i > k$ to get n nonnegative singular values: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{n-1} \geq \sigma_n = 0$. And you can extend $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ with any orthonormal basis $\{\mathbf{v}_{k+1}, \dots, \mathbf{v}_n\}$ for the nullspace of A to get a full set of right singular vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_k, \mathbf{v}_{k+1}, \dots, \mathbf{v}_n\}$.

One last piece of terminology: We will say a k -dimensional subspace V is **best-fit** or **optimal** for a matrix A if it minimizes the sum of the squares of the perpendicular distances between the rows of A and all k -dimensional subspaces. Equivalently, a subspace V is **best-fit** or **optimal** for A if it maximizes the sum of the squared lengths of the projections of the rows of A onto all k -dimensional subspaces.

OK, finally, let's state and prove the "alternate" form of SVD:

Theorem: Alternate Singular Value Decomposition

Let $A_{m \times n} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix}$, $\mathbf{x}_i \in \mathbb{R}^n$, where the rank of A is r . Let $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ be its right singular vectors and let $\{\sigma_1, \dots, \sigma_r\}$ be its singular values:

$$\mathbf{v}_i, \sigma_i \equiv \operatorname{argmax}, \max_{\mathbf{v} \perp \mathbf{v}_1, \dots, \mathbf{v} \perp \mathbf{v}_{i-1}, \|\mathbf{v}\|=1} \|A\mathbf{v}\|$$

For any $k \leq r$, $V_k \equiv \operatorname{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is the best-fit k -dimensional vector subspace for A .

SVD Proof: One At A Time

For $k = 1$, we have already shown that this reduces to a least squares fitting problem for which \mathbf{v}_1 is the solution. For $k = 2$, let W be a best-fit 2-dimensional subspace for A . For any basis $\mathbf{w}_1, \mathbf{w}_2$ of W , $\|A\mathbf{w}_1\| + \|A\mathbf{w}_2\|$ is the sum of squared lengths of the projections of the rows of A onto W . Now,

choose a basis $\mathbf{w}_1, \mathbf{w}_2$ of W so that \mathbf{w}_2 is perpendicular to \mathbf{v}_1 . If \mathbf{v}_1 is perpendicular to W , any unit vector in W will do as \mathbf{w}_2 . If not, choose \mathbf{w}_2 to be the unit vector in W perpendicular to the projection of \mathbf{v}_1 onto W . Since \mathbf{v}_1 was chosen to maximize $\|A\mathbf{v}\|^2$, it follows that $\|A\mathbf{w}_1\|^2 \leq \|A\mathbf{v}_1\|^2$. Since \mathbf{v}_2 was chosen to maximize $\|A\mathbf{v}\|^2$ over all \mathbf{v} perpendicular to \mathbf{v}_1 , $\|A\mathbf{w}_2\|^2 \leq \|A\mathbf{v}_2\|^2$. Hence

$$\|A\mathbf{w}_1\|^2 + \|A\mathbf{w}_2\|^2 \leq \|A\mathbf{v}_1\|^2 + \|A\mathbf{v}_2\|^2$$

and V_2 is at least as good as W and so is a best-fit 2-dimensional subspace.

For general k , proceed by induction. By the induction hypothesis, V_{k-1} is a best-fit $k-1$ dimensional subspace. Suppose W is a best-fit k -dimensional subspace. Choose a basis $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k$ of W so that \mathbf{w}_k is perpendicular to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$. Then

$$\|A\mathbf{w}_1\|^2 + \|A\mathbf{w}_2\|^2 + \dots + \|A\mathbf{w}_k\|^2 \leq \|A\mathbf{v}_1\|^2 + \|A\mathbf{v}_2\|^2 + \dots + \|A\mathbf{v}_{k-1}\|^2 + \|A\mathbf{w}_k\|^2$$

since V_{k-1} is an optimal $k-1$ dimensional subspace. Since \mathbf{w}_k is perpendicular to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$, by the definition of \mathbf{v}_k , $\|A\mathbf{w}_k\|^2 \leq \|A\mathbf{v}_k\|^2$. Thus

$$\|A\mathbf{w}_1\|^2 + \|A\mathbf{w}_2\|^2 + \dots + \|A\mathbf{w}_{k-1}\|^2 + \|A\mathbf{w}_k\|^2 \leq \|A\mathbf{v}_1\|^2 + \|A\mathbf{v}_2\|^2 + \dots + \|A\mathbf{v}_{k-1}\|^2 + \|A\mathbf{v}_k\|^2$$

proving that V_k is at least as good as W and hence is optimal. ■

Now we would like to find a reasonably efficient algorithm for computing the singular vectors and singular values.

The Power Method

For (just about) any $\mathbf{x} \in \mathbb{R}^n$ and any $m \times n$ matrix A :

$$\lim_{s \rightarrow \infty} \frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|} = \mathbf{v}_1 \quad \text{if and only if } \sigma_1 > \sigma_2$$

and

$$\lim_{s \rightarrow \infty} \frac{(A A^T)^s \mathbf{x}}{\|(A A^T)^s \mathbf{x}\|} = \mathbf{u}_1 \quad \text{if and only if } \sigma_1 > \sigma_2$$

where \mathbf{v}_1 is the top right singular vector, \mathbf{u}_1 is the top left singular vector, and σ_1 and σ_2 are the first and second singular values of A .

Proof

First we need to find a convenient formula for $(A^T A)^s$. Suppose A has rank r and orthonormal right singular vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$. Then

$$\begin{aligned} (A^T A)^2 &= \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \right)^2 = \left(\sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^r \sigma_j^2 \mathbf{v}_j \mathbf{v}_j^T \right) \\ &= \sum_{i,j=1}^r \sigma_i^2 \sigma_j^2 \mathbf{v}_i \mathbf{v}_i^T \mathbf{v}_j \mathbf{v}_j^T = \sum_{i=1}^r \sigma_i^4 \mathbf{v}_i \mathbf{v}_i^T \end{aligned}$$

Inductively, we find that

$$(A^T A)^s = \sum_{i=1}^r \sigma_i^{2s} \mathbf{v}_i \mathbf{v}_i^T$$

Aside: The outer product is more convenient here but we also could have found

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$$

The last equation follows from the orthonormality of U so that $U^T U = I$. Also Σ is diagonal and hence symmetric: $\Sigma^T = \Sigma$. Similarly, we can compute

$$(A^T A)^2 = V \Sigma^2 V^T V \Sigma^2 V^T = V \Sigma^4 V^T$$

And inductively find

$$(A^T A)^s = V \Sigma^{2s} V^T$$

Anyway, back to the proof: since $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is an orthonormal basis for \mathbb{R}^n , there exists scalars c_1, \dots, c_n such that $\mathbf{x} = \sum_{i=1}^n c_i \mathbf{v}_i$. And we can compute

$$A^T A \mathbf{x} = \left(\sum_{i=1}^n \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^n c_j \mathbf{v}_j \right) = \sum_{i,j=1}^n \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T c_j \mathbf{v}_j = \sum_{i=1}^n c_i \sigma_i^2 \mathbf{v}_i = \sum_{i=1}^r c_i \sigma_i^2 \mathbf{v}_i$$

Similarly, using our handy formula for $(A^T A)^s$, we find that

$$(A^T A)^s \mathbf{x} = \sum_{i=1}^r c_i \sigma_i^{2s} \mathbf{v}_i$$

We also need a handy formula for $\|(A^T A)^s \mathbf{x}\|$:

$$\begin{aligned} \|(A^T A)^s \mathbf{x}\|^2 &= ((A^T A)^s \mathbf{x})^T ((A^T A)^s \mathbf{x}) = \left(\sum_{i=1}^r c_i \sigma_i^{2s} \mathbf{v}_i \right)^T \left(\sum_{j=1}^r c_j \sigma_j^{2s} \mathbf{v}_j \right) \\ &= \left(\sum_{i=1}^r c_i \sigma_i^{2s} \mathbf{v}_i^T \right) \left(\sum_{j=1}^r c_j \sigma_j^{2s} \mathbf{v}_j \right) = \sum_{i,j=1}^r c_i c_j \sigma_i^{2s} \sigma_j^{2s} \mathbf{v}_i^T \mathbf{v}_j = \sum_{i=1}^r c_i^2 \sigma_i^{4s} \end{aligned}$$

So we have a tractable formula for $\frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|}$ as well:

$$\begin{aligned} \frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|} &= \frac{1}{\|(A^T A)^s \mathbf{x}\|} \sum_{i=1}^r c_i \sigma_i^{2s} \mathbf{v}_i = \frac{1}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} \sum_{i=1}^r c_i \sigma_i^{2s} \mathbf{v}_i \\ &= \sum_{i=1}^r \frac{c_i \sigma_i^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} \mathbf{v}_i \end{aligned}$$

So let's compute $\lim_{s \rightarrow \infty} \frac{c_i \sigma_i^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}}$:

$$\lim_{s \rightarrow \infty} \frac{c_i \sigma_i^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} = \lim_{s \rightarrow \infty} \sqrt{\frac{c_i^2 \sigma_i^{4s}}{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} = \sqrt{\lim_{s \rightarrow \infty} \frac{c_i^2 \sigma_i^{4s}}{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}}$$

$$\begin{aligned}
&= \sqrt{\lim_{s \rightarrow \infty} \frac{1}{1 + \frac{1}{c_i^2 \sigma_i^{4s}} \sum_{j=1, j \neq i}^r c_j^2 \sigma_j^{4s}}} = \sqrt{\lim_{s \rightarrow \infty} \frac{1}{1 + \sum_{j=1, j \neq i}^r \frac{c_j^2 \sigma_j^{4s}}{c_i^2 \sigma_i^{4s}}}} \\
&= \sqrt{\lim_{s \rightarrow \infty} \frac{1}{1 + \sum_{j=1, j \neq i}^r \left(\frac{c_j}{c_i}\right)^2 \left(\frac{\sigma_j}{\sigma_i}\right)^{4s}}}
\end{aligned}$$

For $i = 1$, we have

$$\lim_{s \rightarrow \infty} \frac{c_1 \sigma_1^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} = \begin{cases} 1 & \text{if } \sigma_1 > \sigma_2 \\ \sqrt{\frac{1}{1 + (\frac{c_2}{c_1})^2}} & \text{if } \sigma_1 = \sigma_2 \end{cases}$$

Of course we are also assuming that the random vector \mathbf{x} has a nonzero first component $c_1 \neq 0$ (this is the "just about" disclaimer in the theorem statement). For $i = 2$, we have

$$\lim_{s \rightarrow \infty} \frac{c_2 \sigma_2^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} = \begin{cases} 0 & \text{if } \sigma_1 > \sigma_2 \\ \sqrt{\frac{1}{1 + (\frac{c_1}{c_2})^2}} & \text{if } \sigma_1 = \sigma_2 \end{cases}$$

And putting it together:

$$\lim_{s \rightarrow \infty} \frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|} = \lim_{s \rightarrow \infty} \sum_{i=1}^r \frac{c_i \sigma_i^{2s}}{\sqrt{\sum_{j=1}^r c_j^2 \sigma_j^{4s}}} \mathbf{v}_i = \mathbf{v}_1 \quad \text{so long as } \sigma_1 > \sigma_2$$

The proof for $\lim_{s \rightarrow \infty} \frac{(A A^T)^s \mathbf{x}}{\|(A A^T)^s \mathbf{x}\|} = \mathbf{u}_1$ if and only if $\sigma_1 > \sigma_2$ is similar. ■

Of course, once we have numerically computed \mathbf{v}_1 , we can use the Rayleigh Quotient to get the corresponding eigenvalue of $A^T A$ which is also the square of the singular value of A : $\sigma_1^2 = \lambda_1 = r(\mathbf{v}_1) = \frac{\mathbf{v}_1^T A^T A \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} = \mathbf{v}_1^T A^T A \mathbf{v}_1 = (A \mathbf{v}_1)^T A \mathbf{v}_1 = \mathbf{u}_1^T \mathbf{u}_1$. And then we can normalize \mathbf{u}_1 to finish off the first pair of singular vectors.

Since we have computed σ_1 , \mathbf{v}_1 , and \mathbf{u}_1 , we can repeat this procedure for $A - \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$. The top singular value and vectors of this matrix are σ_2 , \mathbf{v}_2 , and \mathbf{u}_2 . Repeating this procedure k times will get us the top k singular values and vectors and the best-fit k -dimensional subspace $V_k = \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$.

So the question becomes, how quickly does $\frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|}$ converge to \mathbf{v}_1 ? Well, the Power Method isn't industrial strength, but it's still pretty good. You can find a [detailed analysis here](#) or [at the end this](#). It turns out that the convergence is geometric with ratio $\frac{\lambda_2}{\lambda_1}$.

When we actually code this up, when we're computing an approximation for some \mathbf{v}_i , instead of iterating a calculated number of times, we will simply stop iterating when the angle between iteration results is sufficiently small. What does that mean? Let $\mathbf{t}_s = \frac{(A^T A)^s \mathbf{x}}{\|(A^T A)^s \mathbf{x}\|}$. Then \mathbf{t}_s is a unit vector. So the \cos of the angle θ_s between \mathbf{t}_s and \mathbf{t}_{s-1} can be computed as $\cos(\theta_s) = \frac{\mathbf{t}_s \cdot \mathbf{t}_{s-1}}{\|\mathbf{t}_s\| \|\mathbf{t}_{s-1}\|} = \mathbf{t}_s \cdot \mathbf{t}_{s-1}$. We know that $\lim_{\theta \rightarrow 0} \cos(\theta) = 1$ and that $\cos(\theta) \leq 1$ for all θ . At each iteration s , we check whether $\mathbf{t}_s \cdot \mathbf{t}_{s-1} > 1 - \epsilon$ for some arbitrarily small ϵ . When this check is true for some iteration s , we have a sufficiently small angle between iteration results and we're not going to get a much better approximation for \mathbf{v}_i than \mathbf{t}_s . So we stop.