

Chapter 1

SVD, PCA & Pre- processing

Part 1: Linear algebra and SVD



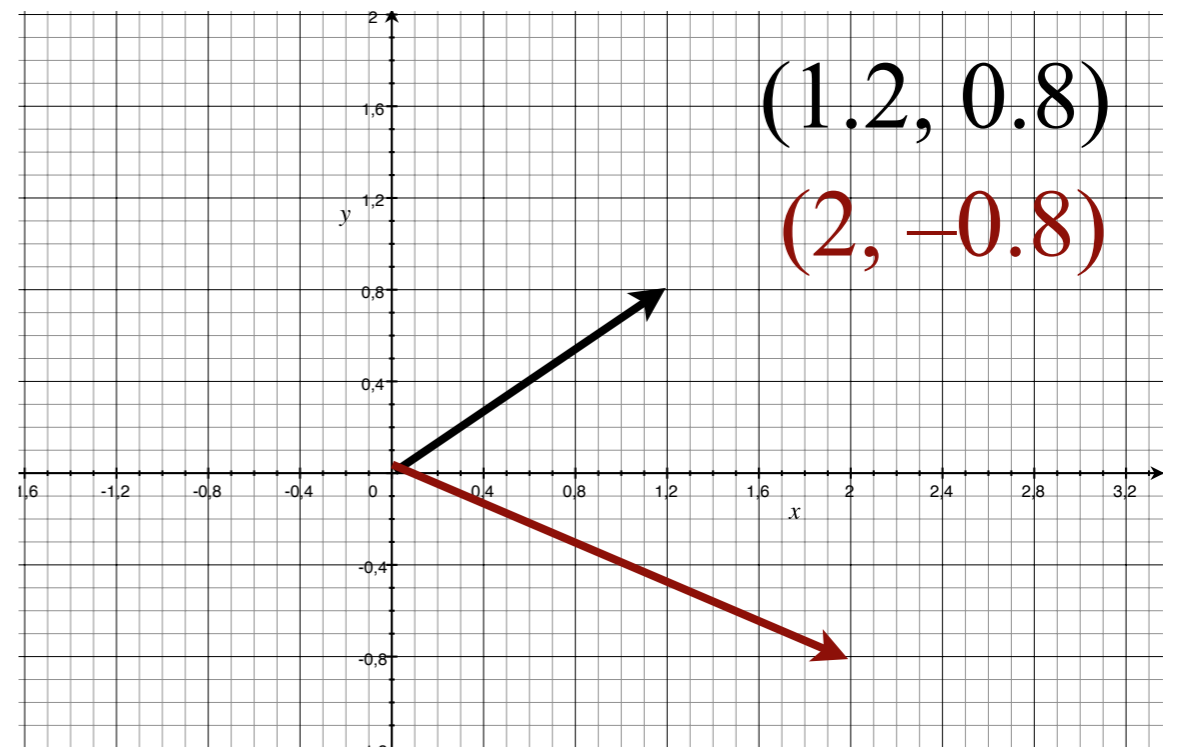
Contents

- Linear algebra crash course
- The singular value decomposition
- Normalization
- Selecting the rank
- The principal component analysis

Linear Algebra Crash Course

Matrices and vectors

- A **vector** is
 - a 1D array of numbers
 - a geometric entity with magnitude and direction
 - a matrix with exactly one row or column



Norms and angles

- The magnitude is measured by a (vector) **norm**

- The **Euclidean** norm

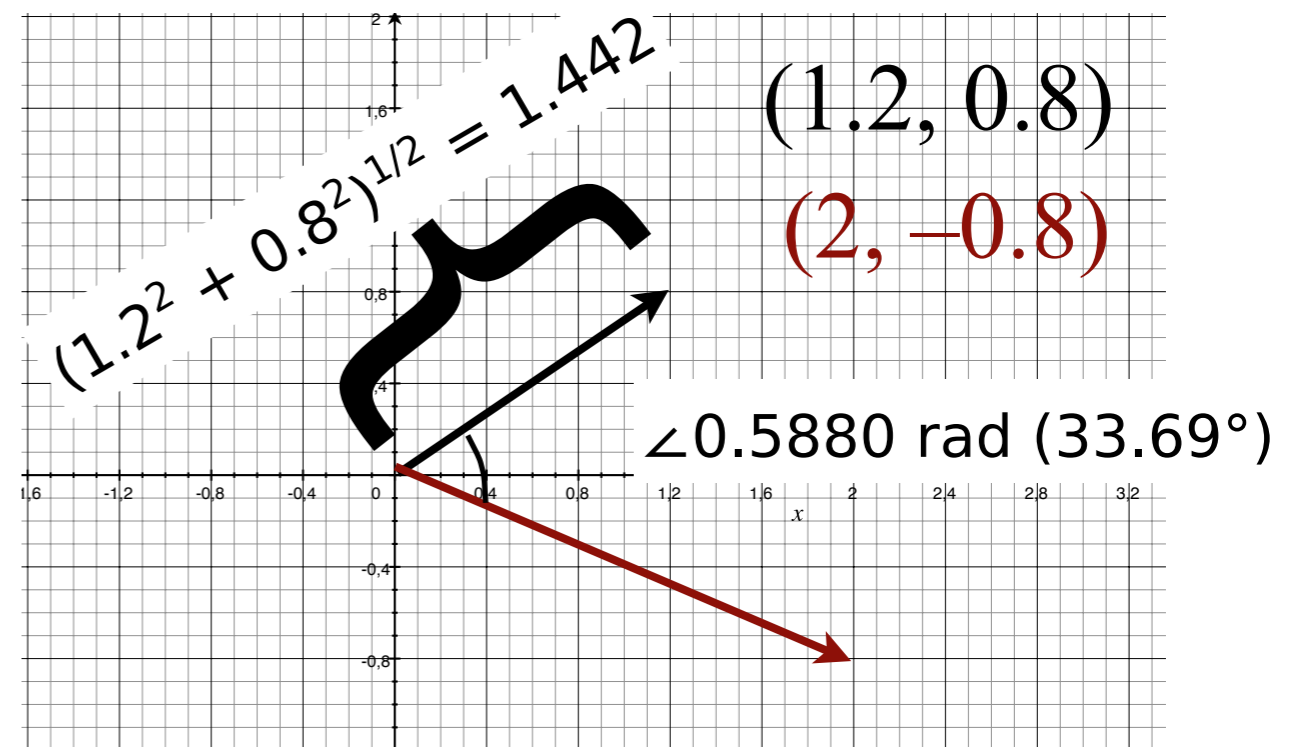
$$\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \left(\sum_{i=1}^n x_i^2\right)^{1/2}$$

- General L_p norm

$$(1 \leq p \leq \infty)$$

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$$

- The direction is measured by the **angle**



Basic vector operations

- The **transpose** of \mathbf{x} , \mathbf{x}^T , transposes a row vector into a column vector and vice versa
- A **dot product** of two vectors of the same dimension is $\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^n x_i y_i$
 - A.k.a. **scalar product** or **inner product**
 - Same as $\langle \mathbf{x}, \mathbf{y} \rangle$, $\mathbf{a}^T \mathbf{b}$ (for column vectors), or $\mathbf{a} \mathbf{b}^T$ (for row vectors)

Orthogonality

- **Orthogonality** is a generalization of perpendicularity
- \mathbf{x} and \mathbf{y} are orthogonal if $\mathbf{x} \cdot \mathbf{y} = 0$
- in Euclidean space: $\mathbf{x} \cdot \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta$
 - θ is the angle between \mathbf{x} and \mathbf{y}

Matrix algebra

- Matrices in $\mathbb{R}^{n \times n}$ form a ring
 - Addition, subtraction, and multiplication
 - But usually no division
 - Multiplication is not commutative
 - **$AB \neq BA$** in general

Matrix multiplication

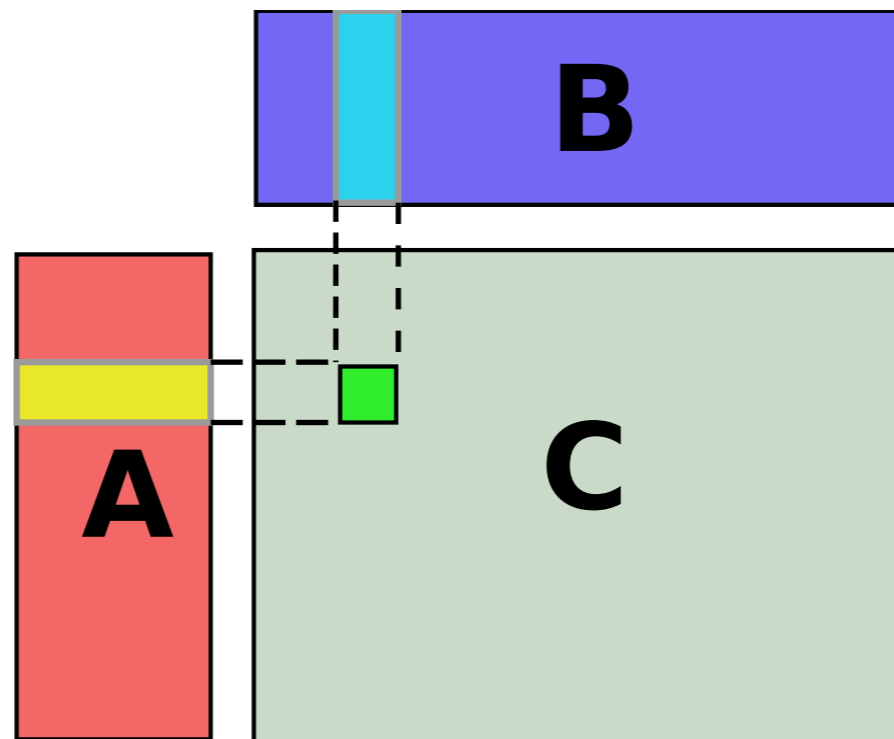
- The product of two matrices, **A** and **B**, is defined element-wise as

$$(\mathbf{AB})_{ij} = \sum_{\ell=1}^k a_{i\ell} b_{\ell j}$$

- The number of columns in **A** and number of rows in **B** must agree
 - inner dimension

Intuition for Matrix Multiplication

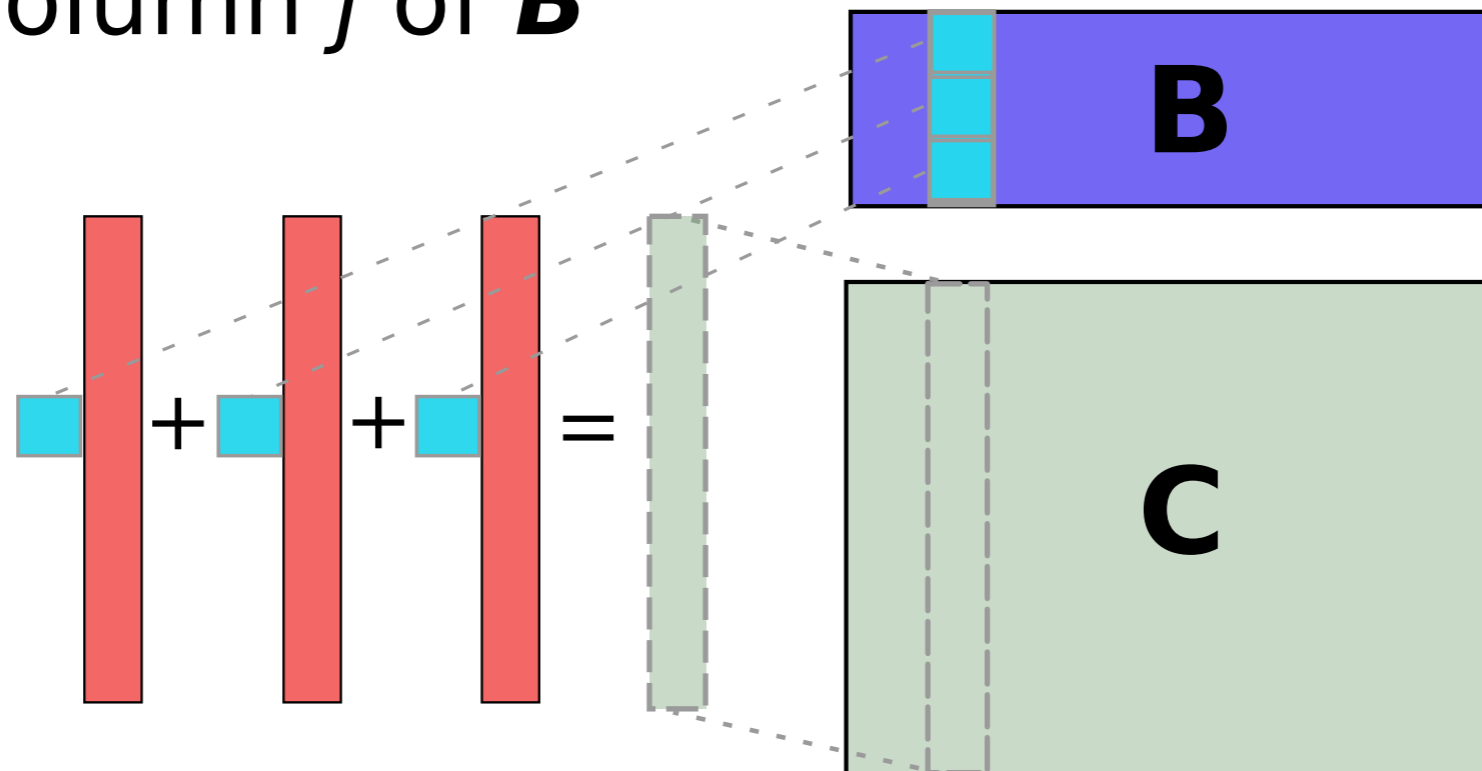
- Element $(\mathbf{AB})_{ij}$ is the inner product of row i of \mathbf{A} and column j of \mathbf{B}



$$c_{ij} = \sum_{\ell=1}^k a_{i\ell} b_{\ell j}$$

Intuition for Matrix Multiplication

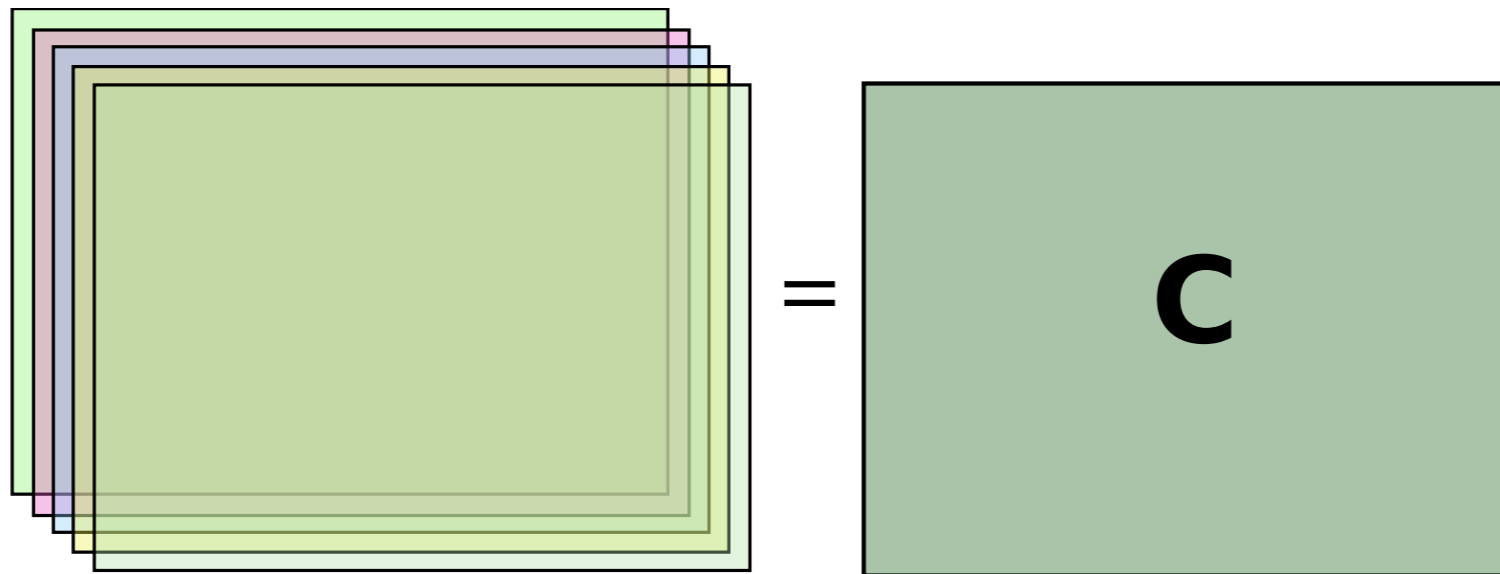
- Column j of \mathbf{AB} is the linear combination of columns of \mathbf{A} with the coefficients coming from column j of \mathbf{B}



$$\mathbf{C} = \left[\left[\sum_{\ell=1}^k b_{\ell 1} \mathbf{a}_{\ell} \right] \quad \left[\sum_{\ell=1}^k b_{\ell 2} \mathbf{a}_{\ell} \right] \cdots \left[\sum_{\ell=1}^k b_{\ell m} \mathbf{a}_{\ell} \right] \right]$$

Intuition for Matrix Multiplication

- Matrix \mathbf{AB} is a sum of k matrices $\mathbf{a}_l \mathbf{b}_l^T$ obtained by multiplying the l -th column of \mathbf{A} with the l -th row of \mathbf{B}



$$\mathbf{C} = \sum_{\ell=1}^k \mathbf{a}_\ell \mathbf{b}_\ell^T$$

Matrix decompositions

- A **decomposition** of matrix **A** expresses it as a product of two (or more) **factor matrices**
 - **$A = BC$**
- Every matrix has decomposition **$A = AI$** (or **$A = IA$** if $n < m$)
- The size of the decomposition is the inner dimension of the product

Matrices as linear maps

- Matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ is a **linear mapping** from \mathbb{R}^m to \mathbb{R}^n
 - $\mathbf{A}(\mathbf{x}) = \mathbf{Ax}$
- If $\mathbf{A} \in \mathbb{R}^{n \times k}$ and $\mathbf{B} \in \mathbb{R}^{k \times m}$, then \mathbf{AB} is a mapping from \mathbb{R}^m to \mathbb{R}^n
- The transpose \mathbf{A}^T is a mapping from \mathbb{R}^n to \mathbb{R}^m
 - $(\mathbf{A}^T)_{ij} = \mathbf{A}_{ji}$
 - $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$

Matrix inverse

- Square matrix **A** is **invertible** if there is a matrix **B** s.t. **AB = BA = I**
 - **B** is the inverse of **A**, denoted **A**⁻¹
 - Usually the transpose is **not** the inverse
- Non-square matrices don't have general inverses
 - Can have left or right inverse:
AR = I or **LA = I**

Linear independency

- Vector \mathbf{u} is **linearly dependent** on a set of vectors $\mathbf{V} = \{\mathbf{v}_i\}$ if \mathbf{u} is a linear combination of \mathbf{v}_i
 - $\mathbf{u} = \sum_i a_i \mathbf{v}_i$ for some a_i
 - If \mathbf{u} is not linearly dependent, it is **linearly independent**
- Set V of vectors is **linearly independent** if all \mathbf{v}_i are linearly independent of $V \setminus \{\mathbf{v}_i\}$

Matrix ranks

- The **column rank** of a matrix \mathbf{A} is the number of linearly independent columns of \mathbf{A}
- The **row rank** of \mathbf{A} is the number of linearly independent rows of \mathbf{A}
- The **Schein rank** of \mathbf{A} is the least integer k such that \mathbf{A} can be expressed as a sum of k rank-1 matrices
 - Rank-1 matrix is an outer product of two vectors

Orthogonal matrices

- Set of vectors $\{\mathbf{v}_i\}$ is **orthogonal** if all \mathbf{v}_i are mutually orthogonal, i.e. $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$ for all $i \neq j$
 - If $\|\mathbf{v}_i\|_2 = 1$ for all \mathbf{v}_i , the set is **orthonormal**
- Square matrix \mathbf{A} is orthogonal if its columns form a set of orthonormal vectors
 - Non-square matrices can be row- or column-orthogonal
- If \mathbf{A} is orthogonal, then $\mathbf{A}^{-1} = \mathbf{A}^T$

Properties of orthogonal matrices

- The inverse of orthogonal matrices is easy to compute
- Orthogonal matrices perform a rotation
 - Only the angle of the vector is changed, the length stays the same

Matrix norms

- **Matrix norms** measure the magnitude of the matrix
 - the magnitude of the values or the image

- **Operator norms:**

$$\|\mathbf{A}\|_p = \max\{\|\mathbf{M}\mathbf{x}\|_p : \|\mathbf{x}\|_p = 1\} \text{ for } p \geq 1$$

- **Frobenius norm:**

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

Singular Value Decomposition

“The SVD is the Swiss Army knife of matrix decompositions”

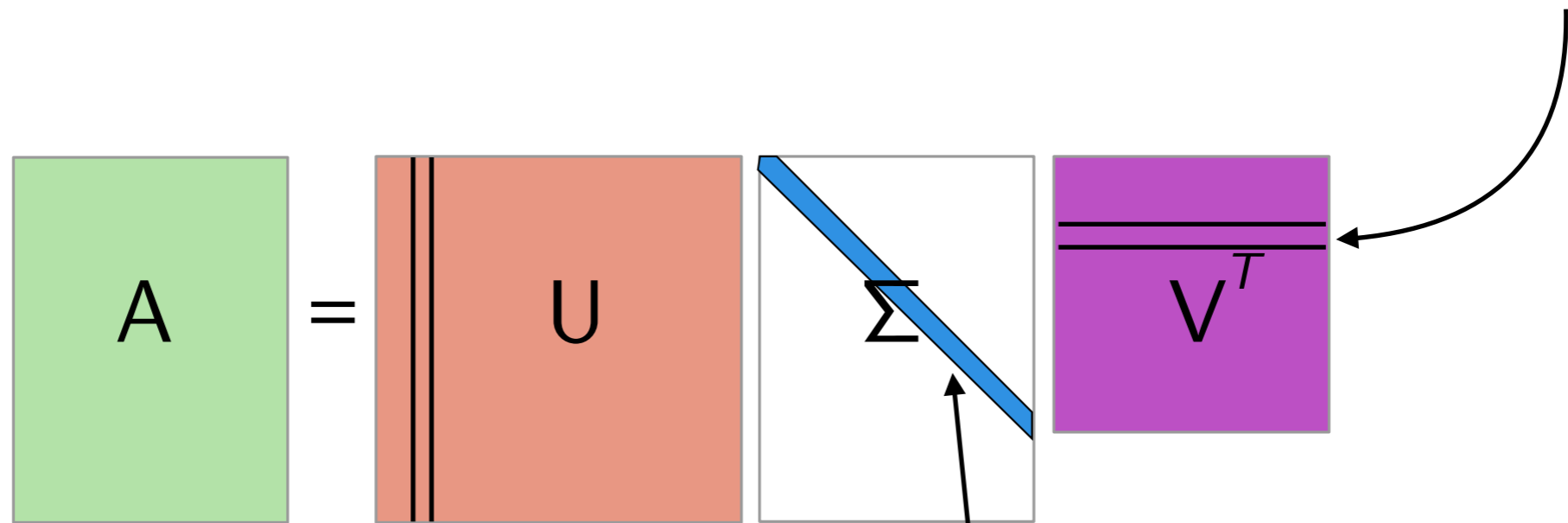
– Diane O’Leary, 2006

The definition

- **Theorem.** For every $\mathbf{A} \in \mathbb{R}^{n \times m}$ there exists an n -by- n orthogonal matrix \mathbf{U} and an m -by- m orthogonal matrix \mathbf{V} such that $\mathbf{U}^T \mathbf{A} \mathbf{V}$ is an n -by- m diagonal matrix $\mathbf{\Sigma}$ that has values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{n,m\}} \geq 0$ in its diagonal
- I.e. every \mathbf{A} has decomposition $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$
- The **singular value decomposition** of \mathbf{A}

In picture

\mathbf{v}_i are the **right singular vectors**



σ_i are the **singular values**

\mathbf{u}_i are the **left singular vectors**

Some useful equations

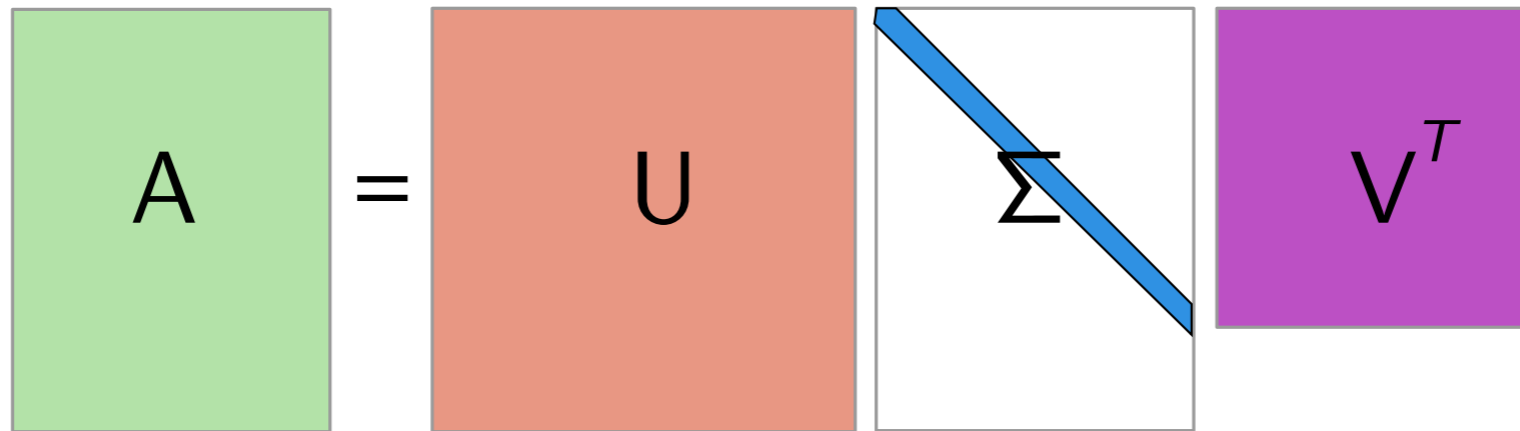
- $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$
 - Expresses \mathbf{A} as a sum of rank-1 matrices
- $\mathbf{A}^{-1} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$ (if \mathbf{A} is invertible)
- $\mathbf{A}^T \mathbf{A} \mathbf{v}_i = \sigma_i^2 \mathbf{v}_i$ (for any \mathbf{A})
- $\mathbf{A} \mathbf{A}^T \mathbf{u}_i = \sigma_i^2 \mathbf{u}_i$ (for any \mathbf{A})

Truncated SVD

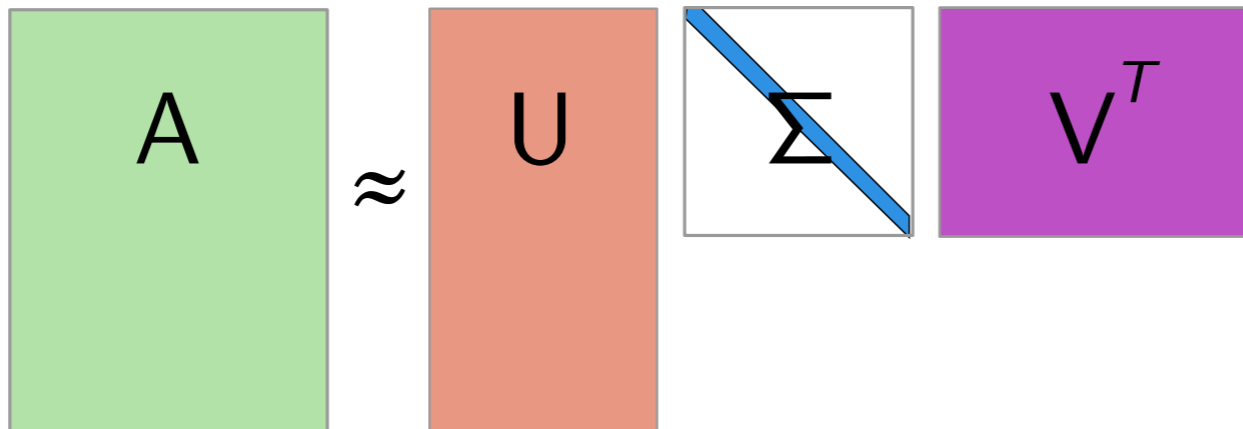
- The rank of the matrix is the number of its non-zero singular values (write $\mathbf{A} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$)
- The **truncated SVD** takes the first k columns of \mathbf{U} and \mathbf{V} and the main k -by- k submatrix of $\mathbf{\Sigma}$
 - $\mathbf{A}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T$
 - \mathbf{U}_k and \mathbf{V}_k are column-orthogonal

Truncated SVD

Full



Truncated



Why is SVD important?

- It gives us the **dimensions of the fundamental subspaces**
- It lets us **compute various norms**
- It tells about **sensitivity of linear systems**
- It gives us optimal solutions to **least-squares linear systems**
- It gives us the **least-error rank- k decomposition**
- **Every matrix has one**

Fundamental theorem of linear algebra

- **Theorem.** Every n -by- m matrix \mathbf{A} induces four **fundamental subspaces**
 - The **range** of dimension $\text{rank}(\mathbf{A}) = r$
 - The set of all linear combinations of columns of \mathbf{A}
 - The **kernel** of dimension $m - r$
 - The set of all vectors \mathbf{x} for which $\mathbf{Ax} = \mathbf{0}$
 - The **coimage** of dimension r
 - The **cokernel** of dimension $n - r$

Fundamental subspaces

- The bases for the fundamental subspaces are:
 - Range: the first r columns of \mathbf{U}
 - Kernel: the last $(m - r)$ columns of \mathbf{V}
 - Coimage: the first r columns of \mathbf{V}
 - Cokernel: the last $(n - r)$ columns of \mathbf{U}

SVD and norms

- Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the SVD of \mathbf{A} .
 - $\|\mathbf{A}\|_F^2 = \sum_{i=1}^{\min\{n,m\}} \sigma_i^2$
 - $\|\mathbf{A}\|_2 = \sigma_1$
- Therefore $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{\min\{n,m\}} \|\mathbf{A}\|_2$
- For truncated SVD, $\|\mathbf{A}_k\|_F^2 = \sum_{i=1}^k \sigma_i^2$

Sensitivity of linear systems

- The solution for system $\mathbf{Ax} = \mathbf{b}$ is $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$
 - Requires that \mathbf{A} is invertible
- Hence $\mathbf{x} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1}\mathbf{b} = \sum_{i=1}^n \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i$
 - Small changes in \mathbf{A} or \mathbf{b} yield large changes in \mathbf{x} if σ_n is small
 - Can we characterize this sensitivity?

Condition number

- The **condition number** $\kappa_p(\mathbf{A})$ of a square matrix \mathbf{A} is $\|\mathbf{A}\|_p \|\mathbf{A}^{-1}\|_p$
 - Particularly $\kappa_2(\mathbf{A}) = \sigma_1(\mathbf{A})/\sigma_n(\mathbf{A})$
 - $\kappa_2(\mathbf{A}) = \infty$ for singular \mathbf{A}
- If κ is large, the matrix is **ill-conditioned**
 - The solution is sensible for small perturbations

Least-squares linear systems

- **Problem.** Given $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^n$, find $\mathbf{x} \in \mathbb{R}^m$ minimizing $\|\mathbf{Ax} - \mathbf{b}\|_2$.
- If \mathbf{A} is invertible, $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is an exact solution
- For non-invertible \mathbf{A} we have to find other solution

The Moore–Penrose pseudo-inverse

- n -by- m matrix \mathbf{B} is the **Moore–Penrose pseudo-inverse** of n -by- m matrix \mathbf{A} if
 - $\mathbf{ABA} = \mathbf{A}$ (but possibly $\mathbf{AB} \neq \mathbf{I}$)
 - $\mathbf{BAB} = \mathbf{B}$
 - $(\mathbf{AB})^T = \mathbf{AB}$ (\mathbf{AB} is symmetric)
 - $(\mathbf{BA})^T = \mathbf{BA}$
- Pseudo-inverse of \mathbf{A} is denoted by \mathbf{A}^+

Pseudo-inverse and SVD

- If $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ is the SVD of \mathbf{A} , then
$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$$
 - $\mathbf{\Sigma}^{-1}$ replaces non-zero σ_i 's with $1/\sigma_i$ and transposes the result
 - N.B. not a real inverse
- **Theorem.** Setting $\mathbf{x} = \mathbf{A}^+\mathbf{y}$ gives the optimal solution to $\|\mathbf{Ax} - \mathbf{y}\|$

The Eckart–Young theorem

- **Theorem.** Let $\mathbf{A}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T$ be the rank- k truncated SVD of \mathbf{A} . Then \mathbf{A}_k is the closest rank- k matrix of \mathbf{A} in the Frobenius sense, that is,

$$\|\mathbf{A} - \mathbf{A}_k\|_F \leq \|\mathbf{A} - \mathbf{B}\|_F \text{ for all rank-}k \text{ matrices } \mathbf{B}$$

- Holds for any unitarily invariant norm

That's all for today

- Next week: normalization and selecting the rank
 - Lecture starts at 12:00 sharp
 - Will end earlier as well
- But SVD will return...