

Linear Algebra Basics

Mahdi Roozbahani

(Covered by) Nimisha Roy

Lecturer, SCI, College of Computing, Georgia Tech

Director, Online Undergraduate Initiatives

Outline

- Linear Algebra Basics ←
- Norms
- Multiplications
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition
- Matrix Calculus

Why linear algebra?

- Most data can be represented or stored in matrix-vector form
- Provides compact representation for sets of linear equations

$$\begin{aligned} 4x_1 - 5x_2 &= -13 \\ -2x_1 + 3x_2 &= 9 \end{aligned} \rightarrow \begin{bmatrix} 4 & -5 \\ -2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -13 \\ 9 \end{bmatrix} \rightarrow \mathbf{Ax} = \mathbf{b}$$

- $\mathbf{A} \in \mathbb{R}^{N \times D}$ denotes a matrix with N rows and D columns, where elements belong to real numbers
- $\mathbf{x} \in \mathbb{R}^D$ denotes a vector with D real entries.

Linear algebra basics

- Transpose of a matrix results from flipping the rows and the columns. Given $\mathbf{A} \in \mathbb{R}^{N \times D}$, the transpose is $\mathbf{A}^T \in \mathbb{R}^{D \times N}$

$$\mathbf{A} = \begin{bmatrix} 4 & -5 \\ -2 & 3 \end{bmatrix} \rightarrow \mathbf{A}^T = \begin{bmatrix} 4 & -2 \\ -5 & 3 \end{bmatrix}$$


- For each element of the matrix, the transpose can be written as $A_{ij}^T = A_{ji}$
- The following properties of the transposes are easily verified
 - $(\mathbf{A}^T)^T = \mathbf{A}$
 - $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$
 - $(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$

Linear algebra basics

- A square matrix $\mathbf{A} \in \mathbb{R}^{D \times D}$ is symmetric if $\mathbf{A} = \mathbf{A}^T$ and it is skew-symmetric if $\mathbf{A} = -\mathbf{A}^T$. Thus each matrix can be written as a sum of symmetric and anti-symmetric matrices:

$$\mathbf{A} = \frac{1}{2}(\mathbf{A} - \mathbf{A}^T) + \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)$$

Outline

- Linear Algebra Basics
- Norms 
- Multiplications
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition

Norms

- Norm of a vector $\|\mathbf{x}\|$ is informally a measure of the length of a vector
- More formally, a norm is any function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ that satisfies:
 - For all $\mathbf{x} \in \mathbb{R}^D$, $f(\mathbf{x}) \geq 0$ (non-negativity)
 - $f(\mathbf{x}) = 0$ if and only if $\mathbf{x} = \mathbf{0}$ (definiteness)
 - For $\mathbf{x} \in \mathbb{R}^D$, $t \in \mathbb{R}$, $f(t\mathbf{x}) = |t|f(\mathbf{x})$ (homogeneity)
 - For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$, $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (triangle inequality)

Norms

- Common norms used in machine learning are:

- ℓ_2 -norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^D x_i^2}$

- ℓ_1 -norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^D |x_i|$

- ℓ_∞ -norm: $\|\mathbf{x}\|_\infty = \max_i |x_i|$

Norms

- All norms presented so far are examples of the family of ℓ_p norms, which are parametrized by a real number $p \geq 1$

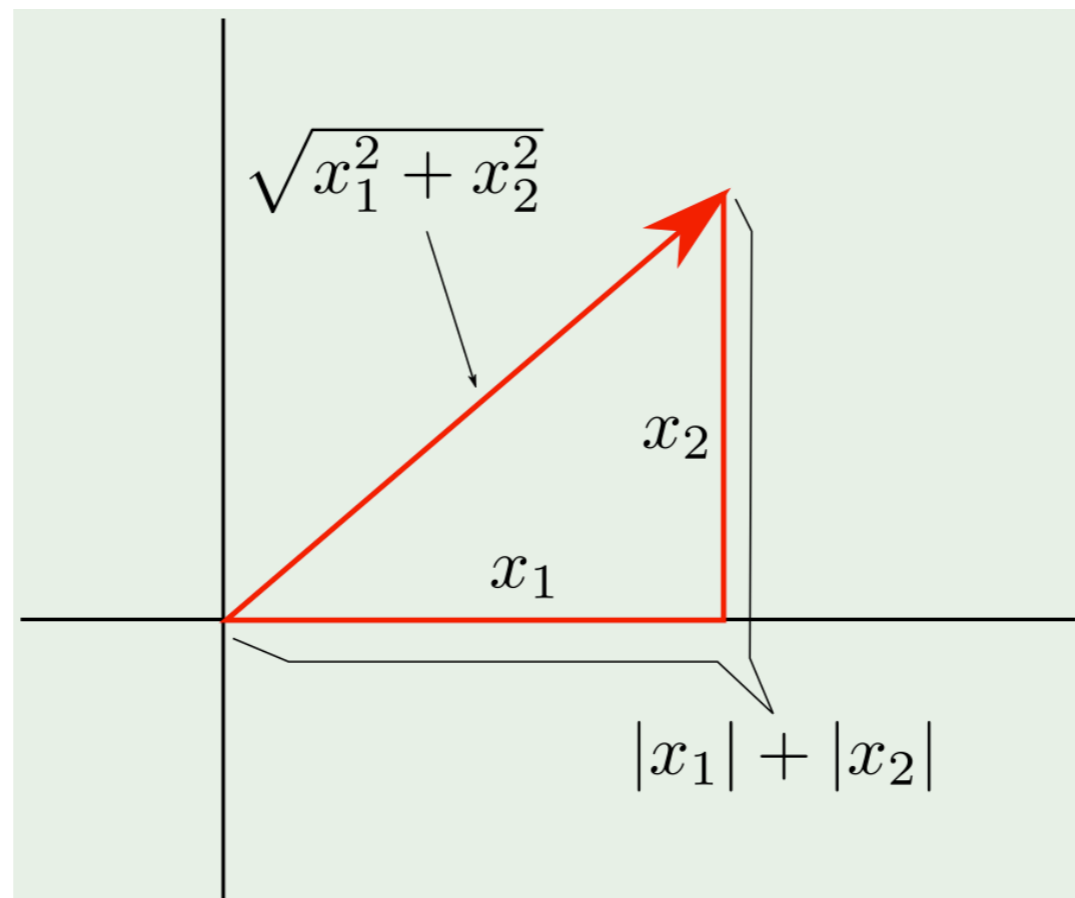
- ℓ_p -norm: $\|\mathbf{x}\|_p = \left(\sum_{i=1}^D |x_i|^p\right)^{\frac{1}{p}}$

- Norms can be defined for matrices, such as the Frobenius norm

- $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^N \sum_{j=1}^D A_{ij}^2} = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})}$

Vector norm examples

- ℓ_1 -norm and ℓ_2 -norm



In context: how are norms useful in ML?

Unsupervised learning

In context: choosing an appropriate norm

Manhattan vs. Euclidean

Special Matrices

- The identity matrix, denoted by $I \in \mathbb{R}^{d \times d}$ is a square matrix with ones on the diagonal and zeros everywhere else
- A diagonal matrix is a matrix where all non-diagonal 'ELEMENTS' are 0. This is typically denoted as $D = \text{diag}(d_1, d_2, \dots, d_d)$

Why is diagonal matrix helpful?


- $A^T A =$ elementwise squaring of all diagonal elements
- $A^{-1} = 1/A$ (elementwise)

Special Matrices

- Two vectors $x, y \in \mathbb{R}^d$ are orthogonal if $x \cdot y = 0$. A square matrix $U \in \mathbb{R}^{d \times d}$ is **Orthonormal** if all its columns are orthogonal to each other and are normalized
- It follows from orthogonality and normality that
 - $U^T U = I = U U^T$
 - $\|Ux\|_2 = \|x\|_2$

Is the inverse of a unitary matrix equal to its transpose?

Outline

- Linear Algebra Basics
- Norms
- Multiplications 
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition

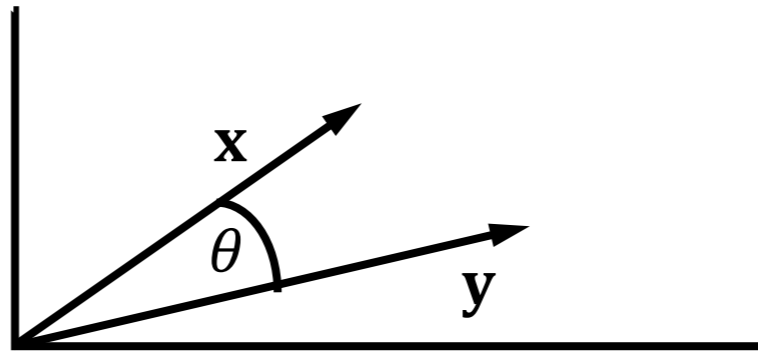
Multiplications

- The product of two matrices $\mathbf{A} \in \mathbb{R}^{N \times D}$ and $\mathbf{B} \in \mathbb{R}^{D \times P}$ is given by $\mathbf{C} \in \mathbb{R}^{N \times P}$, where $C_{ij} = \sum_{k=1}^D A_{ik} B_{kj}$
- Given two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$, the term $\mathbf{x} \cdot \mathbf{y}$ (or $\mathbf{x}^T \mathbf{y}$) is called the inner product or dot product of the vectors, and is a real number given by $\sum_{k=1}^D x_k y_k$. For example,

$$\mathbf{x}^T \mathbf{y} = [x_1 \quad x_2 \quad x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \sum_{i=1}^3 x_i y_i$$

Multiplications

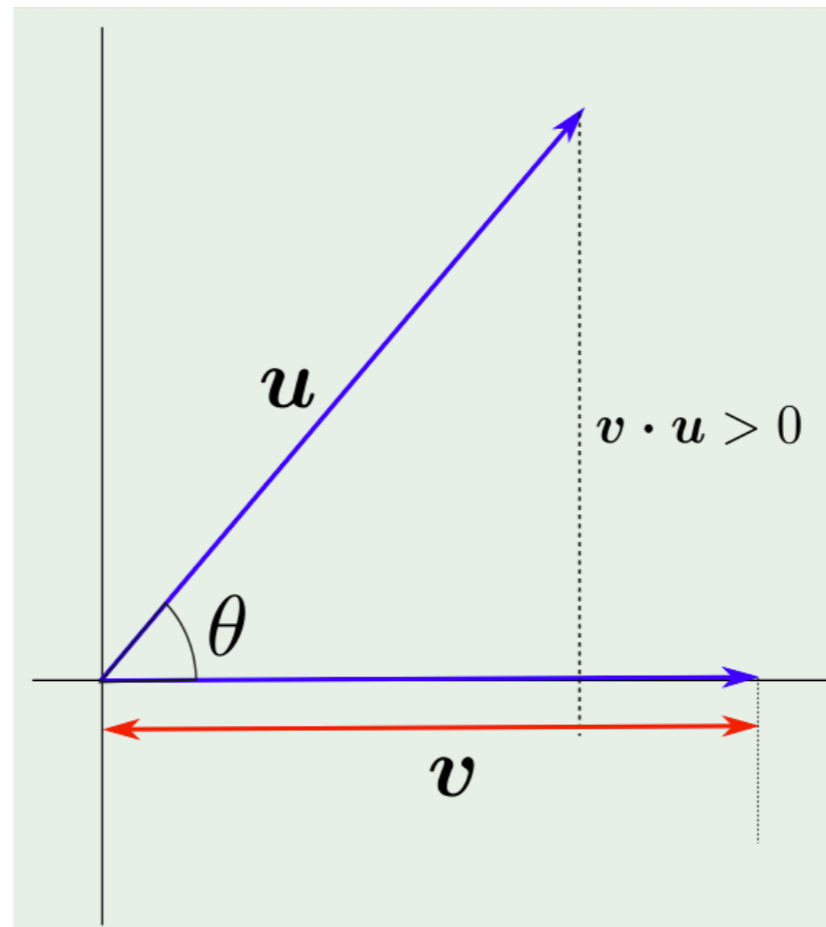
- The dot product also has a geometrical interpretation, for vectors in $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ with an angle θ between them:



$$\mathbf{x} \cdot \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta$$

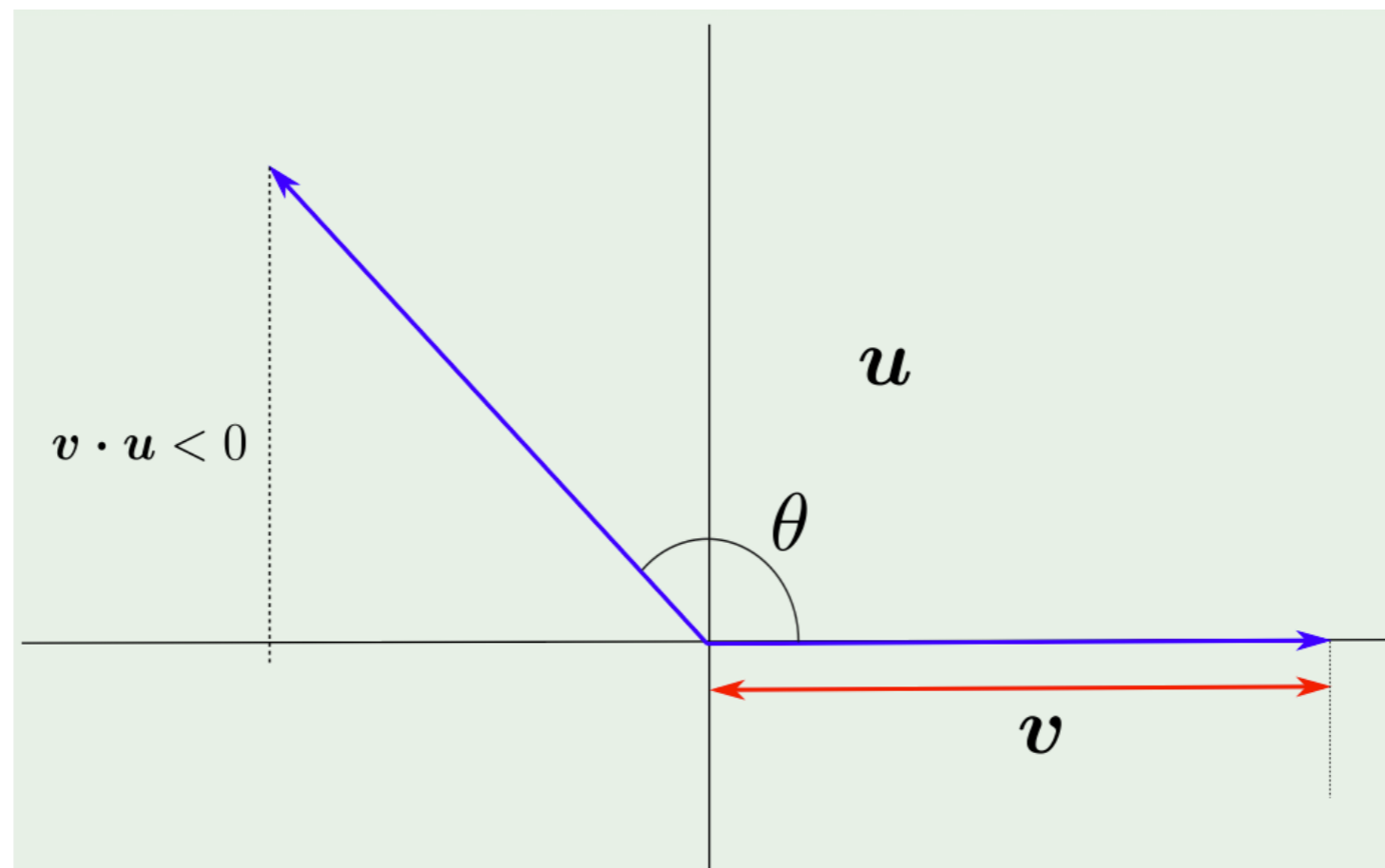
Inner product properties

- The inner product is a measure of **correlation** between two vectors scaled by the norms of the vectors
- Here **correlation** term is used in the loose sense of directional alignment – not statistical sense



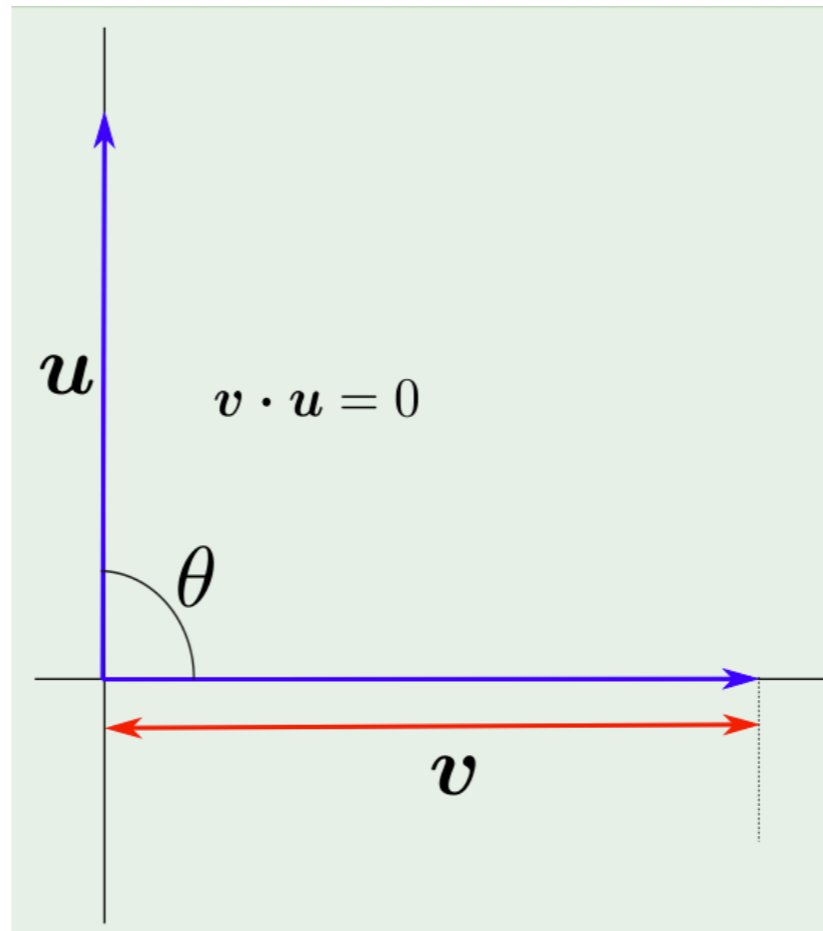
Inner product properties

- The inner product is a measure of correlation between two vectors, scaled by the norms of the vectors



Inner product properties

- The inner product is a measure of correlation between two vectors scaled by the norms of the vectors



**In context: how is the inner product useful
in ML?**

Projecting data onto a new direction


This will be helpful in dimensionality reduction, classification and feature engineering

Multiplications

- Given two vectors $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^D$ the term $\mathbf{xy}^T \in \mathbb{R}^{N \times D}$ is called the outer product of the vectors and is a matrix given by $(x_i y_j)^T = x_i y_j$.
- For example, $\mathbf{x} \in \mathbb{R}^3$ and $\mathbf{y} \in \mathbb{R}^2$

$$\mathbf{xy}^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [y_1 \quad y_2] = \begin{bmatrix} x_1 y_1 & x_1 y_2 \\ x_2 y_1 & x_2 y_2 \\ x_3 y_1 & x_3 y_2 \end{bmatrix}$$

Outline

- Linear Algebra Basics
- Norms
- Multiplications
- Matrix Inversion 
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Matrix Decomposition

Linear Independence and Matrix Rank

- A set of vectors $\{x_1, x_2, \dots, x_d\} \subset \mathbb{R}^d$ are said to be **(linearly) independent** if no vector can be represented as a linear combination of the remaining vectors. That is if

$$x_d = \sum_{i=1}^{d-1} \alpha_i x_i$$

for some scalar values $\alpha_1, \alpha_2, \dots \in \mathbb{R}$ then we say that the vectors are linearly **dependent**; otherwise the vectors are linearly independent

In ML, we want as many linearly independent columns as possible

Linear Independence and Matrix Rank

- The **column rank** of a matrix $A \in \mathbb{R}^{n \times d}$ is the size of the largest subset of columns of A that constitute a linearly independent set. **Row rank** of a matrix is defined similarly for rows of a matrix.

It can be easily shown that the row and column ranks are equivalent, therefore we shall refer only to the **rank** of a matrix.

In general, for a full rank rectangular matrix, rank is the min of number of rows and number of columns.

Matrix Rank: Examples

What are the ranks for the following matrices? How about an identity matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} 1 & 0 & 2 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix}$$

Rank Deficiency means that the number of linearly independent vectors in our set is smaller than the smallest dimension


Matrix Inverse

- The inverse of a square matrix $A \in \mathbb{R}^{d \times d}$ is denoted A^{-1} and is the unique matrix such that $A^{-1}A = I = AA^{-1}$
- For some square matrices A^{-1} may not exist, and we say that A is **singular or non-invertible**. In order for A to have an inverse, A must be **full rank**.
- For non-square matrices the inverse, denoted by A^+ , is given by $A^+ = (A^T A)^{-1} A^T$ called the **pseudo inverse**

In context: why is matrix rank important in ML?

Dataset preprocessing

Outline

- Linear Algebra Basics
- Norms
- Multiplications
- Matrix Inversion
- Trace and Determinant 
- Eigen Values and Eigen Vectors
- Singular Value Decomposition

Matrix Trace

- The trace of a matrix $A \in \mathbb{R}^{d \times d}$, denoted as $\mathbf{tr}(A)$, is the sum of the diagonal elements in the matrix

$$\mathbf{tr}(A) = \sum_{i=1}^d A_{ii}$$

- The trace has the following properties
 - For $A \in \mathbb{R}^{d \times d}$, $\mathbf{tr}(A) = \mathbf{tr}A^\top$
 - For $A, B \in \mathbb{R}^{d \times d}$, $\mathbf{tr}(A + B) = \mathbf{tr}(A) + \mathbf{tr}(B)$
 - For $A \in \mathbb{R}^{d \times d}$, $t \in \mathbb{R}$, $\mathbf{tr}(tA) = t \cdot \mathbf{tr}(A)$
 - For A, B, C such that ABC is a square matrix $\mathbf{tr}(ABC) = \mathbf{tr}(BCA) = \mathbf{tr}(CAB)$
- The trace of a matrix helps us easily compute norms and eigenvalues of matrices as we will see later

Matrix Determinant

Definition (Determinant)

The determinant of a square matrix A , denoted by $|A|$, is defined as

$$\det(A) = \sum_{j=1}^d (-1)^{i+j} a_{ij} M_{ij}$$

where M_{ij} is determinant of matrix A without the row i and column j .

For a 2×2 matrix $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$

$$|A| = ad - bc$$

*Linearly dependent columns \rightarrow matrix is not full rank \rightarrow
determinant 0 \rightarrow non-invertible matrix*

Properties of Matrix Determinant

Basic Properties

- $|A| = |A^T|$
- $|AB| = |A| |B|$
- $|A| = 0$ if and only if A is not invertible
- If A is invertible, then $|A^{-1}| = \frac{1}{|A|}$.

In context: how is the determinant useful in ML?

Matrix inversion

Testing out if $\det(A) = 0$ before trying to invert a matrix saves on computation

Recap

- Special Matrices
- Norms
- Matrix Multiplications
- Matrix Inversion
- Trace and Determinant

Outline

- Linear Algebra Basics
- Norms
- Multiplications
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition



Eigenvalues and Eigenvectors

- Given a square matrix $A \in \mathbb{R}^{d \times d}$ we say that $\lambda \in \mathbb{C}$ is an eigenvalue of A and $x \in \mathbb{C}^d$ is an eigenvector if

$$Ax = \lambda x, \quad x \neq 0$$

- Intuitively this means that upon multiplying the matrix A with a vector x , we get the same vector, but scaled by a parameter λ
- So, eigenvectors are the special vectors whose direction is preserved under transformation (no rotation)

Matrix Eigen Decomposition

- All the eigenvectors can be written together as $AX = X\Lambda$ where the columns of X are the eigenvectors of A , and Λ is a diagonal matrix whose elements are eigenvalues of A

Matrix Eigen Decomposition

- If the eigenvectors of A are invertible, then $A = X\Lambda X^{-1}$
- There are several properties of eigenvalues and eigenvectors
 - $Tr(A) = \sum_{i=1}^d \lambda_i$
 - $|A| = \prod_{i=1}^d \lambda_i$
 - If A is non-singular then $1/\lambda_i$ are the eigenvalues of A^{-1}
 - The eigenvalues of a diagonal matrix are the diagonal elements of the matrix itself!

Matrix Eigen Decomposition

- For a symmetric matrix \mathbf{A} it can be shown that eigenvalues are real and the eigenvectors are orthonormal. Thus it can be represented as $\mathbf{X}\mathbf{\Lambda}\mathbf{X}^T$

Eigenvalues and Eigenvectors

Geometrically, we are transforming the matrix A (if symmetric) from its original orthonormal basis/coordinates to a new set of orthonormal basis x with magnitude as λ

- If A is symmetric, eigenvectors are orthogonal. So, eigenvectors form an orthonormal basis
- Eigenvectors define the directions along which the transformation acts independently
- In the eigenvector basis, there is no rotation or shear, only scaling by the eigenvalues
- Each eigenvalue λ_i controls the amount of stretching or compression along eigenvector x_i

Computing Eigenvalues and Eigenvectors

- We can rewrite the original equation in the following manner

$$\begin{aligned} Ax &= \lambda x, & x &\neq 0 \\ \Rightarrow (A - \lambda I) x &= 0, & x &\neq 0 \end{aligned}$$

- This is only possible if $(A - \lambda I)$ is singular, that is $|(A - \lambda I)| = 0$.
- Thus, eigenvalues and eigenvectors can be computed.
 - Compute the determinant of $A - \lambda I$.
 - This results in a polynomial of degree d .
 - Find the roots of the polynomial by equating it to zero.
 - The d roots are the d eigenvalues of A . They make $A - \lambda I$ singular.
 - For each eigenvalue λ , solve $(A - \lambda I) x$ to find an eigenvector x

Eigenvalue Example

$$\text{Matrix } \mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & -4 \end{bmatrix}$$

1. Compute the determinant of $\mathbf{A} - \lambda\mathbf{I}$

$$\mathbf{A} - \lambda\mathbf{I} = \begin{bmatrix} 1 & 2 \\ 3 & -4 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} 1 - \lambda & 2 \\ 3 & -4 - \lambda \end{bmatrix}$$

$$|\mathbf{A} - \lambda\mathbf{I}| = (1 - \lambda)(-4 - \lambda) - 6$$

2. Find the roots of the polynomial equating it to zero

$$|\mathbf{A} - \lambda\mathbf{I}| = 0 \rightarrow (1 - \lambda)(-4 - \lambda) - 6 = 0 \rightarrow \begin{cases} \lambda_1 = -5 \\ \lambda_2 = 2 \end{cases}$$

Eigenvalue Example

3. For each eigenvalue λ solve $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$ to find eigenvector \mathbf{x}

$$\begin{bmatrix} 1 - \lambda & 2 \\ 3 & -4 - \lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow \begin{cases} (1 - \lambda)x_1 + 2x_2 = 0 \\ 3x_1 - (4 + \lambda)x_2 = 0 \end{cases}$$

Eigenvector for $\lambda_1 = -5$

$$\begin{cases} 6x_1 + 2x_2 = 0 \\ 3x_1 + x_2 = 0 \end{cases} \rightarrow \mathbf{x}_1 = \begin{bmatrix} 1 \\ -3 \end{bmatrix} \text{ or } \begin{bmatrix} 0.3162 \\ -0.9487 \end{bmatrix}$$

Eigenvector for $\lambda_2 = 2$

$$\begin{cases} -x_1 + 2x_2 = 0 \\ 3x_1 - 6x_2 = 0 \end{cases} \rightarrow \mathbf{x}_2 = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \text{ or } \begin{bmatrix} 0.8944 \\ 0.4472 \end{bmatrix}$$

Can a matrix have the same eigenvalues?

If two vectors are linearly independent, does it mean they are orthogonal to each other?

In context: how is eigenvalues and eigenvectors helpful in ML?

- Eigenvectors are “special directions” where a matrix doesn’t rotate or shear – it only stretches or shrinks.
- The stretch amount is the eigenvalue.
- Many ML tools boil down to understanding **data transformations** (via matrices). Eigenvectors/eigenvalues tell us the **essential structure** of those transformations:
 - Eigenvectors reveal the fundamental directions along which the transformation acts independently.
 - Eigenvalues quantify the strength of that action (how much stretching or shrinking occurs).
 - This decomposition simplifies complex transformations into basic “rotate–stretch–rotate back” operations, making it easier to analyze patterns, reduce dimensionality, or understand stability and sensitivity in models.

Outline

- Linear Algebra Basics
- Norms
- Multiplications
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition



Singular Value Decomposition

$$\bar{X}_{n \times d}$$

n: datapoints

d: dimensions

X is a centered matrix

$$\bar{X} = U \Sigma V^T$$

$U_{n \times n} \rightarrow$ unitary matrix $\rightarrow U \times U^T = I$

$\Sigma_{n \times d} \rightarrow$ diagonal matrix

$V_{d \times d} \rightarrow$ unitary matrix $\rightarrow V \times V^T = I$

$$\begin{array}{c}
 X = \begin{bmatrix} u_{1 \times 1} & \dots & \dots & \dots & u_{1 \times n} \\ \vdots & \ddots & \dots & \dots & \vdots \\ \vdots & \vdots & \ddots & \dots & \vdots \\ \vdots & \vdots & \dots & \ddots & \vdots \\ u_{1 \times 1} & \dots & \dots & \dots & u_{n \times n} \end{bmatrix} \times \begin{bmatrix} \Sigma_{1 \times 1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Sigma_{d \times d} \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} v_{1 \times 1} & \dots & \dots & \dots & v_{1 \times d} \\ \vdots & \ddots & \dots & \dots & \vdots \\ \vdots & \vdots & \ddots & \dots & \vdots \\ \vdots & \vdots & \dots & \ddots & \vdots \\ v_{d \times 1} & \dots & \dots & \dots & v_{d \times d} \end{bmatrix} \\
 \begin{array}{ccc} U & \Sigma & V^T \\ & d < n & \end{array}
 \end{array}$$

Covariance matrix:

$$C_{d \times d} = \frac{\bar{X}^T \bar{X}}{n}$$

$$\left. \begin{aligned} \bar{X} &= U \Sigma V^T \\ C &= \frac{\bar{X}^T \bar{X}}{n} \end{aligned} \right\} C = \frac{V \Sigma^T U^T U \Sigma V^T}{n} = \frac{V \Sigma^2 V^T}{n}$$

$$C = \frac{V\Sigma^2V^T}{n} = V \frac{\Sigma^2}{n} V^T$$

$$CV = V \frac{\Sigma^2}{n} V^T V = V \frac{\Sigma^2}{n}$$

$$CV = V\Lambda$$

Remember:

$$AX = X\Lambda$$

$\lambda_i = \frac{\Sigma_i^2}{n} \rightarrow$ The eigenvalues of covariance matrix

λ_i : Eigenvalue of C or covariance matrix

Σ_i : Singular value of X matrix

So, we can **directly** calculate eigenvalue and eigenvectors of a covariance matrix by having the singular value decomposition of matrix X

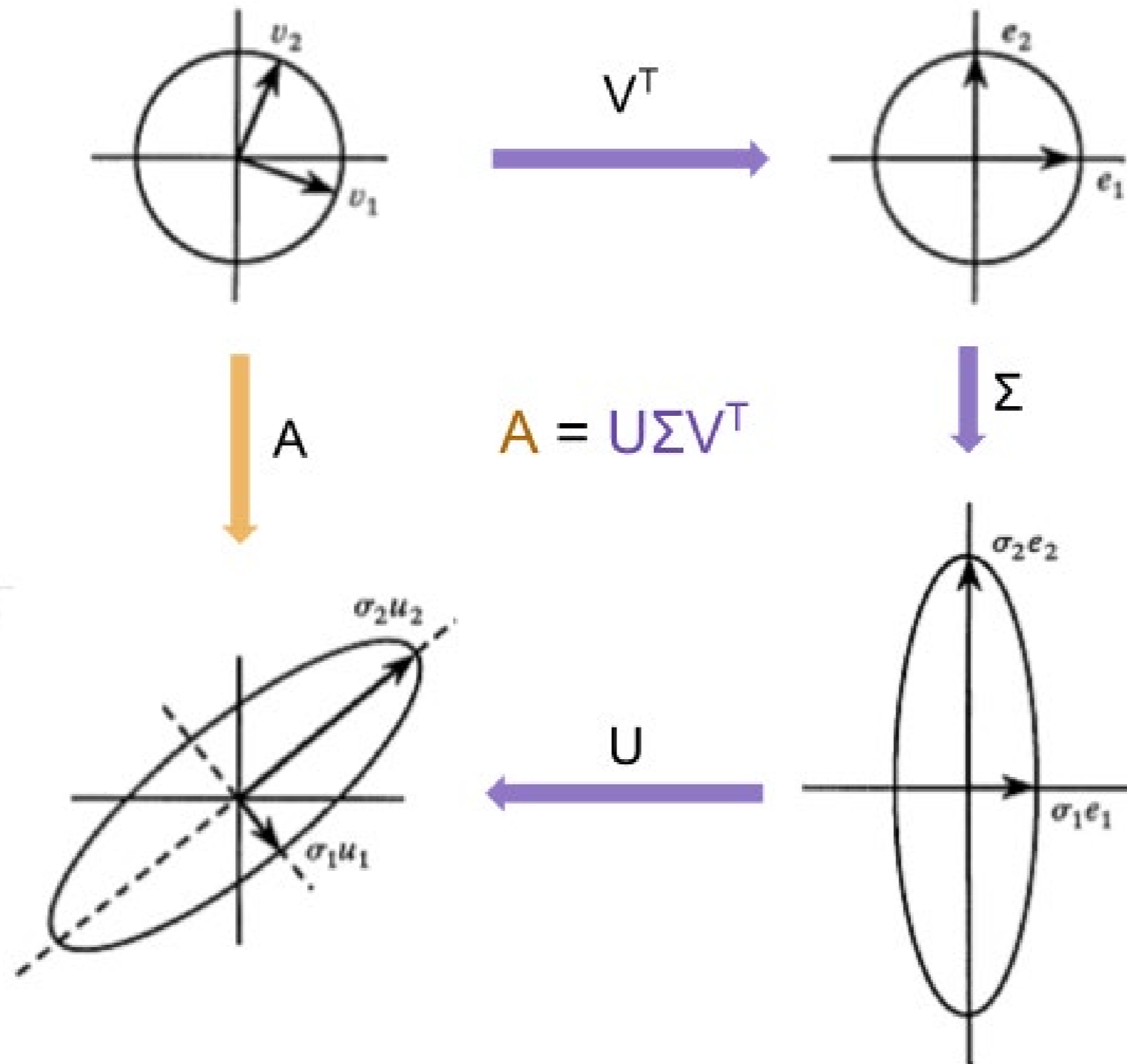
So, we can **directly** calculate eigenvalue and eigenvectors of a covariance matrix by having the singular value decomposition of matrix X

Why is this helpful?

Because SVD **always exists** for **any** matrix (rectangular, non-symmetric) and is **more numerically stable**. That's why in ML (e.g., PCA) we usually prefer **SVD of X** over eigen-decomposition of $X^T X$.

Geometric Meaning of SVD

SVD tells us that *any* linear transformation can be decomposed into:
a rotation \rightarrow a scaling \rightarrow another rotation.



In context: why is SVD and covariance matrix useful in ML?

SVD ($X=U\Sigma V^T$)

- Finds the **most informative directions** in the data (columns of V); σ_i tell how strong each is.
- Gives the **best low-rank approximation** → dimensionality reduction, compression, **denoising** (Low values of σ_i means often represents noise)

Covariance ($C = \frac{X^T X}{n}$ with centered X)

- Summarizes **spread** (variances) and **relationships** (correlations) between features.
- Drives **PCA**: eigenvectors = principal components; eigenvalues = **explained variance**.
- Basis for **whitening/standardization**, which helps many models.

→ Covariance tells us what varies; **SVD** shows how to **use** it—rotate to those directions and keep the big values.

Summary

- Linear Algebra Basics
- Norms
- Multiplications
- Matrix Inversion
- Trace and Determinant
- Eigen Values and Eigen Vectors
- Singular Value Decomposition