

Mathematics of Transformer

Xianbiao Qi

Intellifusion Inc.

November 20, 2024



Mathematical foundations of deep learning

- ① Linear Algebra [22, 1] → Matrix Analysis [10, 30] →
Numerical Linear Algebra [24]
Random Matrix Theory [21]
Matrix Calculus [16]
- ② Calculus → Vector Calculus → Matrix Calculus [16]
- ③ Probability [3] → Stochastic Process → High-dimensional Probability [26]
SDE
- ④ ODE → PDE → SDE [4, 8]
- ⑤ Information Theory [28]
- ⑥ Convex Optimization [14], Numerical Optimization [15]
- ⑦ Mathematics of Machine Learning [6], Probabilistic Machine Learning [13]

Contents

- 1 Part 1: Linear Algebra
[23, 22, 5, 10, 11, 26, 14, 15, 16, 7, 6]
- 2 Part 2: Matrix Calculus
[7, 9, 16, 10, 11]
- 3 Part 3: Transformer
[25, 17, 18, 12, 19, 20, 29, 27, 2]

Part 1: Linear Algebra

[23, 22, 5, 10, 11, 26, 14, 15, 16, 7, 6]

Notation

Definition:

$$\mathbf{x} \in \mathbb{R}^d$$

$$\mathbf{X} \in \mathbb{R}^{d \times n}$$

$$\mathbf{w} \in \mathbb{R}^d$$

$$\mathbf{W} \in \mathbb{R}^{d_1 \times d}$$

Simple cases:

$$y = \mathbf{w}^\top \mathbf{x}, \mathbf{y} = (\mathbf{W}_1 \mathbf{x}), \mathbf{Y} = (\mathbf{W}_q \mathbf{X})^\top (\mathbf{W}_k \mathbf{X})$$

Four key problems in Linear Algebra

- Linear equations

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

- Least square

$$\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

- Eigenvalue decomposition

$$\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$$

- Singular value decomposition (SVD)

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$$

View a matrix in 4 ways

1 Viewing a Matrix – 4 Ways

A matrix ($m \times n$) can be viewed as 1 matrix, mn numbers, n columns and m rows.

$$\begin{bmatrix} \text{gray square} \end{bmatrix} = \begin{bmatrix} \text{blue dots} \\ \text{blue dots} \\ \text{blue dots} \\ \text{blue dots} \end{bmatrix} = \begin{bmatrix} \text{green vertical bars} \\ \text{green vertical bars} \end{bmatrix} = \begin{bmatrix} \text{pink horizontal bars} \\ \text{pink horizontal bars} \\ \text{pink horizontal bars} \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

1 matrix

6 numbers

2 column vectors
with 3 numbers

3 row vectors
with 2 numbers

Figure 1: Viewing a Matrix in 4 Ways

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbf{a}_1 & \mathbf{a}_2 \\ | & | \end{bmatrix} = \begin{bmatrix} -\mathbf{a}_1^* \\ -\mathbf{a}_2^* \\ -\mathbf{a}_3^* \end{bmatrix}$$

View vector-by-vector

v1  =  = ● Dot product (number)

v2  =  =  Rank 1 Matrix

Dot product ($a \cdot b$) is expressed as $a^T b$ in matrix language and yields a number.

ab^T is a matrix ($ab^T = A$). If neither a, b are 0, the result A is a rank 1 matrix.

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = x_1 + 2x_2 + 3x_3$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [x \quad y] = \begin{bmatrix} x & y \\ 2x & 2y \\ 3x & 3y \end{bmatrix}$$

Figure 2: Vector times Vector - (v1), (v2)

View matrix-by-vector

Mv1

$$\begin{bmatrix} \text{pink} \\ \text{pink} \\ \text{pink} \end{bmatrix} \begin{bmatrix} \text{green} \end{bmatrix} = \begin{bmatrix} \text{pink} \\ \text{pink} \\ \text{pink} \end{bmatrix}$$

Mv2

$$\begin{bmatrix} \text{green} \\ \text{green} \end{bmatrix} \begin{bmatrix} \text{blue} \\ \text{blue} \end{bmatrix} = \bullet \begin{bmatrix} \text{green} \end{bmatrix} + \bullet \begin{bmatrix} \text{green} \end{bmatrix}$$

The row vectors of A are multiplied by a vector \mathbf{x} and become the three dot-product elements of $A\mathbf{x}$.

The product $A\mathbf{x}$ is a linear combination of the column vectors of A .

$$A\mathbf{x} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} (x_1+2x_2) \\ (3x_1+4x_2) \\ (5x_1+6x_2) \end{bmatrix}$$

$$A\mathbf{x} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} + x_2 \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix}$$

Figure 3: Matrix times Vector - (Mv1), (Mv2)

View matrix-by-matrix

MM 1

Every element becomes a dot product of row vector and column vector.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = \begin{bmatrix} (x_1+2x_2) & (y_1+2y_2) \\ (3x_1+4x_2) & (3y_1+4y_2) \\ (5x_1+6x_2) & (5y_1+6y_2) \end{bmatrix}$$

MM 2

Ax and Ay are linear combinations of columns of A .

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = A[x \ y] = [Ax \ Ay]$$

MM 3

The produced rows are linear combinations of rows.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^* X \\ \mathbf{a}_2^* X \\ \mathbf{a}_3^* X \end{bmatrix} X = \begin{bmatrix} \mathbf{a}_1^* X \\ \mathbf{a}_2^* X \\ \mathbf{a}_3^* X \end{bmatrix}$$

MM 4

Multiplication AB is broken down to a sum of rank 1 matrices.

$$\begin{aligned} \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} &= [\mathbf{a}_1 \ \mathbf{a}_2] \begin{bmatrix} \mathbf{b}_1^* \\ \mathbf{b}_2^* \end{bmatrix} = \mathbf{a}_1 \mathbf{b}_1^* + \mathbf{a}_2 \mathbf{b}_2^* \\ &= \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} [b_{11} \ b_{12}] + \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix} [b_{21} \ b_{22}] = \begin{bmatrix} b_{11} & b_{12} \\ 3b_{11} & 3b_{12} \\ 5b_{11} & 5b_{12} \end{bmatrix} + \begin{bmatrix} 2b_{21} & 2b_{22} \\ 4b_{21} & 4b_{22} \\ 6b_{21} & 6b_{22} \end{bmatrix} \end{aligned}$$

Figure 6: Matrix times Matrix - (MM1), (MM2), (MM3), (MM4)

Matrix decomposition¹

$A = CR$		Independent columns in C Row echelon form in R Leads to column rank = row rank
$A = LU$		LU decomposition from Gaussian elimination (Lower triangular)(Upper triangular)
$A = QR$		QR decomposition as Gram-Schmidt orthogonalization Orthogonal Q and triangular R
$S = Q\Lambda Q^T$		Eigenvalue decomposition of a symmetric matrix S Eigenvectors in Q , eigenvalues in Λ
$A = U\Sigma V^T$		Singular value decomposition of all matrices A Singular values in Σ

Table 1: The Five Factorization

¹<https://github.com/kenjihiranabe/The-Art-of-Linear-Algebra>

Matrix world

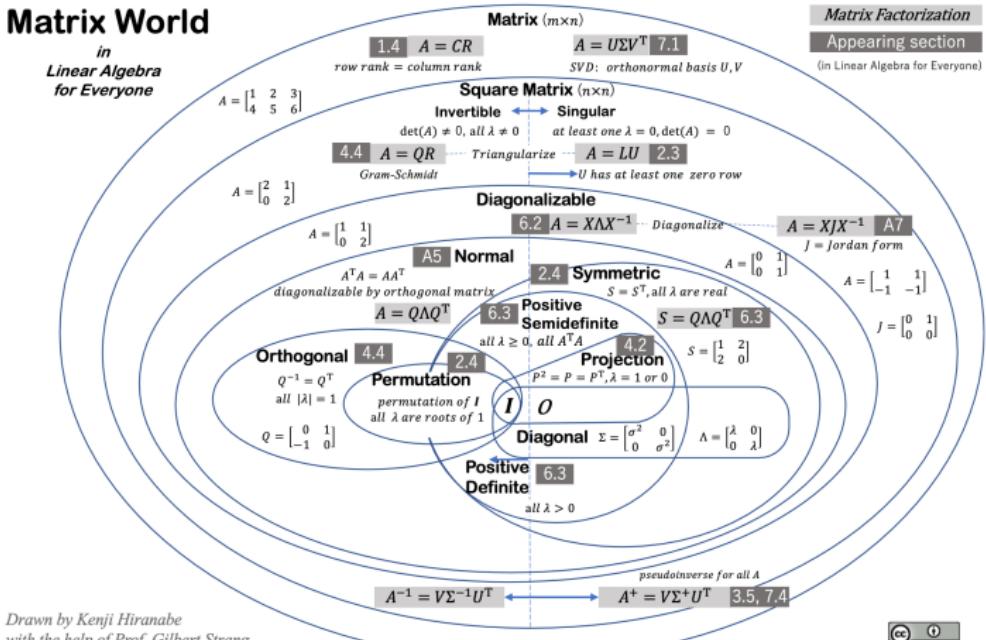


Figure 19: Matrix World

Rotation matrix

Definition

A rotation matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix that preserves distances and orientation:

- $\mathbf{R}^\top \mathbf{R} = \mathbf{R}\mathbf{R}^\top = \mathbf{I}$
- $\det(\mathbf{R}) = 1$
- $\|\mathbf{R}\mathbf{x}\| = \|\mathbf{x}\|$ for all \mathbf{x}

Key Properties

- Inverse equals transpose: $\mathbf{R}^{-1} = \mathbf{R}^\top$
- Composition of rotations is a rotation
- All eigenvalues have magnitude 1
- Preserves angles between vectors

2D rotation matrix

Counterclockwise Rotation by Angle θ

$$\mathbf{R}(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

Action on point (x, y) :

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

Example: 45-degree rotation

Rotating by $\theta = \frac{\pi}{4}$

$$\mathbf{R}(\pi/4) = \begin{pmatrix} \cos(\pi/4) & -\sin(\pi/4) \\ \sin(\pi/4) & \cos(\pi/4) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}$$

Example: Rotating point (1, 0):

$$\begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$$

3D rotation matrix

Rotation Around Coordinate Axes

Around x-axis (Roll):

$$R_x(\theta) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta \\ 0 & \sin \theta & \cos \theta \end{pmatrix}$$

Around y-axis (Pitch):

$$R_y(\theta) = \begin{pmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{pmatrix}$$

Around z-axis (Yaw):

$$R_z(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

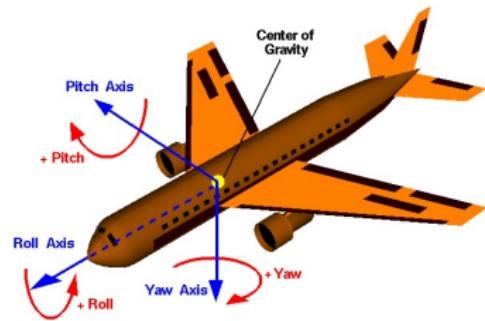


Figure: Roll, Pitch, Yaw.

Determinant

Definition

For a square matrix \mathbf{A} , the determinant $\det(\mathbf{A})$ or $|\mathbf{A}|$ is a scalar value that provides information about the matrix's invertibility and the scaling factor of the linear transformation it represents.

Key Properties

For $n \times n$ matrices \mathbf{A} and \mathbf{B} :

- $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$
- $\det(\mathbf{A}^\top) = \det(\mathbf{A})$
- $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$ if \mathbf{A} is invertible
- For triangular matrices: product of diagonal entries

Computing determinant

For 2×2 matrix:

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

$$\det(\mathbf{A}) = a_{11}a_{22} - a_{12}a_{21}$$

For 3×3 matrix:

$$|\mathbf{A}| = \begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix}$$

Use cofactor expansion

Geometric Interpretation

The determinant represents:

- Area (2D) or Volume (3D) scaling factor

Trace

Definition

For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, the trace of \mathbf{A} , denoted $\text{trace}(\mathbf{A})$, is the sum of its diagonal elements:

$$\text{trace}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$$

Basic Properties

For matrices \mathbf{A}, \mathbf{B} and scalar c :

- Linearity: $\text{trace}(\mathbf{A} + \mathbf{B}) = \text{trace}(\mathbf{A}) + \text{trace}(\mathbf{B})$
- Scalar multiplication: $\text{trace}(c\mathbf{A}) = c \text{trace}(\mathbf{A})$
- Transpose: $\text{trace}(\mathbf{A}) = \text{trace}(\mathbf{A}^\top)$
- Trace of identity: $\text{trace}(\mathbf{I}_n) = n$

Properties of trace

Cyclic Property

For matrices \mathbf{A} and \mathbf{B} of compatible dimensions:

$$\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA})$$

More generally:

$$\text{trace}(\mathbf{ABP}) = \text{trace}(\mathbf{BPA}) = \text{trace}(\mathbf{PAB})$$

Connection to Matrix Operations

- Similar matrices: If $\mathbf{B} = \mathbf{P}^{-1}\mathbf{AP}$, then $\text{trace}(\mathbf{A}) = \text{trace}(\mathbf{B})$
- For outer product: $\text{trace}(\mathbf{xx}^\top) = \mathbf{x}^\top \mathbf{x}$
- Frobenius norm: $\|\mathbf{A}\|_F^2 = \text{trace}(\mathbf{A}^\top \mathbf{A})$

Examples

For 2×2 matrix:

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

$$\text{trace}(\mathbf{A}) = 3 + 4 = 7$$

For block matrix:

$$\mathbf{B} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$$

$$\text{trace}(\mathbf{B}) = \text{trace}(\mathbf{A}_{11}) + \text{trace}(\mathbf{A}_{22})$$

Definition

For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, the rank of \mathbf{A} , denoted $\text{rank}(\mathbf{A})$, can be defined in multiple equivalent ways:

- Dimension of the column space: $\text{rank}(\mathbf{A}) = \dim(\text{col}(\mathbf{A}))$
- Dimension of the row space: $\text{rank}(\mathbf{A}) = \dim(\text{row}(\mathbf{A}))$
- Number of linearly independent columns or rows
- Number of non-zero singular values

Fundamental properties

Basic Properties

For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$:

- $0 \leq \text{rank}(\mathbf{A}) \leq \min(m, n)$
- $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^\top)$
- $\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$
- For square matrices: Full rank \iff invertible

Computing rank

Methods to Find Rank

- Gaussian elimination
- Compute determinants of submatrices
- SVD decomposition: count non-zero singular values

Matrix Decompositions

For a matrix \mathbf{A} with $\text{rank}(\mathbf{A}) = r$:

- Reduced SVD: $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$ with r singular values
- Full rank factorization: $\mathbf{A} = \mathbf{F}\mathbf{G}$ where $\mathbf{F} \in \mathbb{R}^{m \times r}$, $\mathbf{G} \in \mathbb{R}^{r \times n}$

Applications

System of Linear Equations

For the system $\mathbf{A}\mathbf{x} = \mathbf{b}$:

- Solution exists $\iff \text{rank}(\mathbf{A}) = \text{rank}([\mathbf{A}|\mathbf{b}])$
- Unique solution $\iff \text{rank}(\mathbf{A}) = n$ (full column rank)
- No solution $\iff \text{rank}(\mathbf{A}) < \text{rank}([\mathbf{A}|\mathbf{b}])$

Applications

- Linear independence analysis
- Matrix invertibility testing
- Image and signal processing (low-rank approximations)
- Data compression and dimensionality reduction

Examples

Full Rank Matrix:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\text{rank}(\mathbf{A}) = 2$$

Rank Deficient Matrix:

$$\mathbf{B} = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$$

$$\text{rank}(\mathbf{B}) = 1$$

Rank and Matrix Properties

- $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$
- $\text{rank}(\mathbf{A}\mathbf{A}^\top) = \text{rank}(\mathbf{A}^\top\mathbf{A}) = \text{rank}(\mathbf{A})$

Matrix inverse

Definition

For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, its inverse \mathbf{A}^{-1} (if it exists) satisfies:

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

Basic Properties

If \mathbf{A} is invertible:

- $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
- $(\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top$
- $(c\mathbf{A})^{-1} = \frac{1}{c}\mathbf{A}^{-1}$ for scalar $c \neq 0$
- $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ (reverse order!)

Invertibility conditions

Equivalent Conditions

For a square matrix \mathbf{A} , the following are equivalent:

- \mathbf{A} is invertible
- $\det(\mathbf{A}) \neq 0$
- $\text{rank}(\mathbf{A}) = n$
- The columns of \mathbf{A} are linearly independent
- The equation $\mathbf{Ax} = \mathbf{b}$ has a unique solution for all \mathbf{b}
- 0 is not an eigenvalue of \mathbf{A}

Computing inverses

Methods for Finding \mathbf{A}^{-1}

1. Adjugate Method (for 2×2 and 3×3):

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A})$$

2. Gaussian Elimination:

$$[\mathbf{A} | \mathbf{I}] \rightsquigarrow [\mathbf{I} | \mathbf{A}^{-1}]$$

3. For 2×2 matrix:

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \implies \mathbf{A}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

Special inverse matrix

Notable Cases

- Orthogonal matrix: $\mathbf{A}^{-1} = \mathbf{A}^\top$
- Diagonal matrix: \mathbf{A}^{-1} has reciprocal diagonal entries
- Triangular matrix: \mathbf{A}^{-1} has same triangular form
- Symmetric matrix: \mathbf{A}^{-1} is also symmetric

Useful Results

- If \mathbf{A} is invertible: $(\mathbf{A}^k)^{-1} = (\mathbf{A}^{-1})^k$
- Inverse of block diagonal matrix is block diagonal
- For similar matrices: If $\mathbf{B} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}$, then $\mathbf{B}^{-1} = \mathbf{P}^{-1}\mathbf{A}^{-1}\mathbf{P}$

Applications

Solving Linear Systems

For $\mathbf{A}\mathbf{x} = \mathbf{b}$ with invertible \mathbf{A} :

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

Kronecker product

Definition

For $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{p \times q}$, the Kronecker product $\mathbf{A} \otimes \mathbf{B}$ is:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}$$

Dimensions

- Result is $(mp) \times (nq)$ matrix
- Each $a_{ij}\mathbf{B}$ is a $p \times q$ block

Properties

Fundamental Properties

For compatible matrices and scalars:

- Bilinearity: $(c\mathbf{A}) \otimes \mathbf{B} = \mathbf{A} \otimes (c\mathbf{B}) = c(\mathbf{A} \otimes \mathbf{B})$
- Associativity: $(\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C})$
- Not commutative: $\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A}$ in general
- Distributivity: $(\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C}$

Transpose Property

$$(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$$

Important properties

Matrix Operations

For compatible matrices:

- $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$
- $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ (if both invertible)
- $\det(\mathbf{A} \otimes \mathbf{B}) = (\det \mathbf{A})^n (\det \mathbf{B})^m$
- $\text{rank}(\mathbf{A} \otimes \mathbf{B}) = \text{rank}(\mathbf{A}) \text{rank}(\mathbf{B})$

With Identity Matrices

- $\mathbf{I}_m \otimes \mathbf{I}_n = \mathbf{I}_{mn}$
- $(\mathbf{I}_m \otimes \mathbf{A})(\mathbf{I}_m \otimes \mathbf{B}) = \mathbf{I}_m \otimes (\mathbf{AB})$

Vectorization

The Vec Operator

For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{vec}(\mathbf{A})$ stacks the columns of \mathbf{A} into a single vector:

$$\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n] \implies \text{vec}(\mathbf{A}) = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_n \end{bmatrix}$$

where a_i are the columns of \mathbf{A}

Dimensions

- Input: $\mathbf{A} \in \mathbb{R}^{m \times n}$
- Output: $\text{vec}(\mathbf{A}) \in \mathbb{R}^{mn}$

Properties

Linear Properties

For matrices \mathbf{A}, \mathbf{B} and scalar c :

- $\text{vec}(\mathbf{A} + \mathbf{B}) = \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B})$
- $\text{vec}(c\mathbf{A}) = c \text{ vec}(\mathbf{A})$
- $\text{vec}(\mathbf{A}^\top) \neq \text{vec}(\mathbf{A})$ in general

Inner Product Relation

For matrices \mathbf{A}, \mathbf{B} of same size:

$$\text{trace}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$$

Connection of Kronecker product to Vec operator

Vec Operator

For matrix \mathbf{A} , $\text{vec}(\mathbf{A})$ stacks columns of \mathbf{A} into a vector

Important Relations

- $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \mathbf{B}^\top \otimes \mathbf{I}_m) \text{vec}(\mathbf{A}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B}) = (\mathbf{I}_l \otimes \mathbf{AB}) \text{vec}(\mathbf{C})$
- $\text{vec}(\mathbf{AB}) = (\mathbf{I} \otimes \mathbf{A}) \text{vec}(\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{I}) \text{vec}(\mathbf{A})$
- These relations are crucial in matrix equations and optimization

Example

Simple Example

For $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ and $\mathbf{B} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} 1\mathbf{B} & 2\mathbf{B} \\ 3\mathbf{B} & 4\mathbf{B} \end{bmatrix} = \begin{bmatrix} a & b & 2a & 2b \\ c & d & 2c & 2d \\ 3a & 3b & 4a & 4b \\ 3c & 3d & 4c & 4d \end{bmatrix}$$

Kronecker product and Vec operator

Fundamental Identity

For compatible matrices $\mathbf{A}, \mathbf{X}, \mathbf{B}$:

$$\text{vec}(\mathbf{AXB}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X})$$

Special Cases

- $\text{vec}(\mathbf{AX}) = (\mathbf{I} \otimes \mathbf{A}) \text{vec}(\mathbf{X})$
- $\text{vec}(\mathbf{XB}) = (\mathbf{B}^\top \otimes \mathbf{I}) \text{vec}(\mathbf{X})$
- $\text{vec}(\mathbf{Ax}) = (\mathbf{x}^\top \otimes \mathbf{I}) \text{vec}(\mathbf{A})$

Properties of the Vec operator

Key Properties

For matrices \mathbf{A} , \mathbf{B} , \mathbf{X} and vectors \mathbf{a} , \mathbf{c} :

$$\text{vec}(\mathbf{AXB}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X})$$

$$\text{trace}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$$

$$\text{vec}(\mathbf{A} + \mathbf{B}) = \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B})$$

$$\text{vec}(\alpha \mathbf{A}) = \alpha \cdot \text{vec}(\mathbf{A})$$

$$\mathbf{a}^\top \mathbf{XBX}^\top \mathbf{c} = \text{vec}(\mathbf{X})^\top (\mathbf{B} \otimes \mathbf{ca}^\top) \text{vec}(\mathbf{X})$$

Commutation Matrix

Let \mathbf{A} denote the following 3×2 matrix:

$$\mathbf{A} = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}.$$

\mathbf{A} has the following column-major and row-major vectorizations (respectively):

$$\mathbf{v}_{\text{col}} = \text{vec}(\mathbf{A}) = [1, 2, 3, 4, 5, 6]^{\top}, \quad \mathbf{v}_{\text{row}} = \text{vec}(\mathbf{A}^{\top}) = [1, 4, 2, 5, 3, 6]^{\top}$$

The associated commutation matrix is

$$\mathbf{K} = \mathbf{K}^{(3,2)} = \begin{bmatrix} 1 & . & . & . & . & . \\ . & . & . & 1 & . & . \\ . & 1 & . & . & . & . \\ . & . & . & . & 1 & . \\ . & . & 1 & . & . & . \\ . & . & . & . & . & 1 \end{bmatrix},$$

Expectation of linear combinations

Assume \mathbf{X} and \mathbf{x} to be a matrix and a vector of random variables. Then:

$$E[\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C}] = \mathbf{A}E[\mathbf{X}]\mathbf{B} + \mathbf{C}$$

$$\text{Var}[\mathbf{A}\mathbf{x}] = \mathbf{A}\text{Var}[\mathbf{x}]\mathbf{A}^\top$$

$$\text{Cov}[\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}] = \mathbf{A}\text{Cov}[\mathbf{x}, \mathbf{y}]\mathbf{B}^\top$$

Assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , then:

$$E[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbf{m} + \mathbf{b}$$

$$E[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbf{m}$$

$$E[\mathbf{x} + \mathbf{b}] = \mathbf{m} + \mathbf{b}$$

Expectation I

Assume \mathbf{x} to be a stochastic vector with mean \mathbf{m} , and covariance \mathbf{M} .
Then:

$$\mathbb{E}[(\mathbf{A}\mathbf{x} + \mathbf{a})(\mathbf{B}\mathbf{x} + \mathbf{b})^\top] = \mathbf{A}\mathbf{M}\mathbf{B}^\top + (\mathbf{A}\mathbf{m} + \mathbf{a})(\mathbf{B}\mathbf{m} + \mathbf{b})^\top$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{M} + \mathbf{m}\mathbf{m}^\top$$

$$\mathbb{E}[\mathbf{x}\mathbf{a}^\top\mathbf{x}] = (\mathbf{M} + \mathbf{m}\mathbf{m}^\top)\mathbf{a}$$

$$\mathbb{E}[\mathbf{x}^\top\mathbf{a}\mathbf{x}^\top] = \mathbf{a}^\top(\mathbf{M} + \mathbf{m}\mathbf{m}^\top)$$

$$\mathbb{E}[(\mathbf{A}\mathbf{x})(\mathbf{A}\mathbf{x})^\top] = \mathbf{A}(\mathbf{M} + \mathbf{m}\mathbf{m}^\top)\mathbf{A}^\top$$

$$\mathbb{E}[(\mathbf{x} + \mathbf{a})(\mathbf{x} + \mathbf{a})^\top] = \mathbf{M} + (\mathbf{m} + \mathbf{a})(\mathbf{m} + \mathbf{a})^\top$$

Expectation II

Trace formulas for the stochastic vector \mathbf{x} with mean \mathbf{m} and covariance \mathbf{M} :

$$\mathbb{E}[(\mathbf{A}\mathbf{x} + \mathbf{a})^\top(\mathbf{B}\mathbf{x} + \mathbf{b})] = \text{trace}(\mathbf{A}\mathbf{M}\mathbf{B}^\top) + (\mathbf{A}\mathbf{m} + \mathbf{a})^\top(\mathbf{B}\mathbf{m} + \mathbf{b})$$

$$\mathbb{E}[\mathbf{x}^\top \mathbf{x}] = \text{trace}(\mathbf{M}) + \mathbf{m}^\top \mathbf{m}$$

$$\mathbb{E}[\mathbf{x}^\top \mathbf{A}\mathbf{x}] = \text{trace}(\mathbf{A}\mathbf{M}) + \mathbf{m}^\top \mathbf{A}\mathbf{m}$$

$$\mathbb{E}[(\mathbf{A}\mathbf{x})^\top(\mathbf{A}\mathbf{x})] = \text{trace}(\mathbf{A}\mathbf{M}\mathbf{A}^\top) + (\mathbf{A}\mathbf{m})^\top(\mathbf{A}\mathbf{m})$$

$$\mathbb{E}[(\mathbf{x} + \mathbf{a})^\top(\mathbf{x} + \mathbf{a})] = \text{trace}(\mathbf{M}) + (\mathbf{m} + \mathbf{a})^\top(\mathbf{m} + \mathbf{a})$$

Gaussian probability density function

The density of $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$ is

$$p(\mathbf{x}) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp \left[-\frac{1}{2}(\mathbf{x} - \mathbf{m})^\top \Sigma^{-1} (\mathbf{x} - \mathbf{m}) \right]$$

Note that if \mathbf{x} is d -dimensional, then:

$$\det(2\pi\Sigma) = (2\pi)^d \det(\Sigma)$$

Integration and normalization properties follow from this form.

Mean and Covariance

First and second moments. Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$\mathbb{E}(\mathbf{x}) = \mathbf{m}$$

$$\text{Cov}(\mathbf{x}, \mathbf{x}) = \text{Var}(\mathbf{x}) = \Sigma = \mathbb{E}(\mathbf{x}\mathbf{x}^\top) - \mathbb{E}(\mathbf{x})\mathbb{E}(\mathbf{x}^\top) = \mathbb{E}(\mathbf{x}\mathbf{x}^\top) - \mathbf{m}\mathbf{m}^\top$$

As for any other distribution, it holds for gaussians that:

$$\mathbb{E}[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbb{E}[\mathbf{x}]$$

$$\text{Var}[\mathbf{A}\mathbf{x}] = \mathbf{A}\text{Var}[\mathbf{x}]\mathbf{A}^\top$$

$$\text{Cov}[\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}] = \mathbf{A}\text{Cov}[\mathbf{x}, \mathbf{y}]\mathbf{B}^\top$$

1. Linear equations

Assume \mathbf{A} is $m \times n$ and consider the linear system

$$\mathbf{Ax} = \mathbf{b}$$

Construct the augmented matrix $\mathbf{B} = [\mathbf{A}, \mathbf{b}]$,
then

Condition Solution

$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) = n$ Unique solution \mathbf{x}

$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) < n$ Many solutions \mathbf{x}

$\text{rank}(\mathbf{A}) < \text{rank}(\mathbf{B})$ No solutions \mathbf{x}

Basic concepts

Matrix Form

A system of m linear equations in n unknowns:

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

where:

- $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$.

Augmented Matrix

$$[\mathbf{A}|\mathbf{b}] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & | & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & | & b_2 \\ \vdots & \vdots & \ddots & \vdots & | & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} & | & b_m \end{array} \right]$$

Solution existence

Fundamental Theorem

For system $\mathbf{A}\mathbf{x} = \mathbf{b}$:

- Has solution $\iff \text{rank}(\mathbf{A}) = \text{rank}([\mathbf{A}|\mathbf{b}])$
- Has unique solution $\iff \text{rank}(\mathbf{A}) = n$
- Has infinitely many solutions $\iff \text{rank}(\mathbf{A}) < n$ and consistent
- Has no solution $\iff \text{rank}(\mathbf{A}) < \text{rank}([\mathbf{A}|\mathbf{b}])$

Solution methods: direct

Gaussian Elimination

Steps:

- ① Forward elimination (to upper triangular)
- ② Back substitution
- ③ Solution verification

LU Decomposition

$$\mathbf{A} = \mathbf{L}\mathbf{U}$$

- Solve $\mathbf{Ly} = \mathbf{b}$ (forward substitution)
- Solve $\mathbf{Ux} = \mathbf{y}$ (back substitution)
- Efficient for multiple right-hand sides

2. Least square

Objective

Given:

- Data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$
- Target vector $\mathbf{b} \in \mathbb{R}^m$
- Goal: Find $\mathbf{x} \in \mathbb{R}^n$ that minimizes:

$$J(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2 = (\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b})$$

Normal Equations

The optimal solution satisfies:

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}$$

LMS algorithm

Gradient Descent Approach

The gradient of the cost function:

$$\nabla J(x) = 2A^\top(Ax - b)$$

Update rule:

$$x_{k+1} = x_k - \mu \nabla J(x_k)$$

where μ is the learning rate

Algorithm 1 LMS Algorithm

- ```

1: Initialize \mathbf{x}_0 randomly
2: Choose learning rate $\mu > 0$
3: while not converged do
4: $\mathbf{e}_k = \mathbf{A}\mathbf{x}_k - \mathbf{b}$ ▷ Compute error
5: $\mathbf{g}_k = 2\mathbf{A}^\top \mathbf{e}_k$ ▷ Compute gradient
6: $\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \mathbf{g}_k$ ▷ Update estimate
7: end while

```

## Online Learning

Process one sample at a time:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mu(\mathbf{b}_i - \mathbf{a}_i^\top \mathbf{x}_k)\mathbf{a}_i$$

where  $\mathbf{a}_i^\top$  is the  $i$ -th row of  $\mathbf{A}$

## Properties

- Lower memory requirements
- Handles streaming data
- Noisy updates but converges in expectation

### 3. Eigenvalue decomposition

#### Definition

For a square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ :

- A scalar  $\lambda \in \mathbb{C}$  is an eigenvalue if:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

for some nonzero vector  $\mathbf{v} \in \mathbb{C}^n$

- $\mathbf{v}$  is called an eigenvector corresponding to  $\lambda$

#### Characteristic Equation

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0$$

This polynomial equation of degree  $n$  has  $n$  complex roots (counting multiplicity)

# Eigenvalue decomposition

## Definition

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is diagonalizable if:

$$\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$$

where:

- $\mathbf{P}$  is the matrix of eigenvectors
- $\mathbf{D}$  is a diagonal matrix of eigenvalues
- Columns of  $\mathbf{P}$  are eigenvectors of  $\mathbf{A}$

## Key Properties

- $\mathbf{A}^k = \mathbf{P}\mathbf{D}^k\mathbf{P}^{-1}$
- $\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i$
- $\text{trace}(\mathbf{A}) = \sum_{i=1}^n \lambda_i$

## Example: $2 \times 2$ case

For Matrix

$$\mathbf{A} = \begin{pmatrix} 3 & 1 \\ 1 & 3 \end{pmatrix}$$

Eigenvalues:

$$\lambda_1 = 4, \lambda_2 = 2$$

Eigenvectors:

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{v}_2 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

Decomposition:

$$\mathbf{A} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}^{-1}$$

## 4. Singular value decomposition (SVD)

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

where,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ ,  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ .

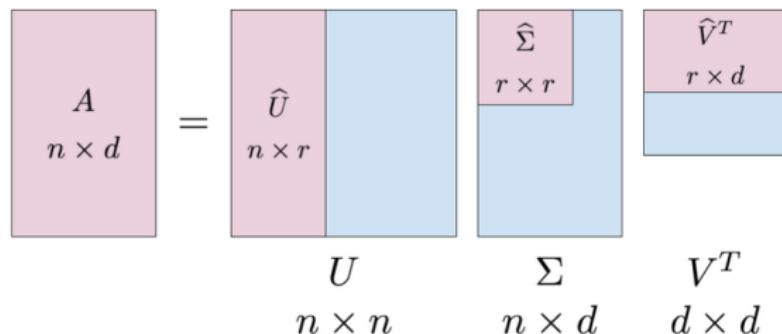


Figure: Illustration<sup>2</sup> of Singular Value Decomposition.

<sup>2</sup><https://www.linkedin.com/pulse/ml-algorithm-singular-value-decomposition-angela-ju/> Navigation icons

## 4. SVD

```
#Demonstrate code of SVD
import torch

Create a simple matrix
X = torch.tensor([[1., 2., 3.],
 [4., 5., 6.]])

Perform SVD
U, S, V = torch.svd(X)

Reconstruct the matrix
X_reconstructed = U @ torch.diag(S) @ V.T

print("Original matrix:")
print(X)
print("\nReconstructed matrix:")
print(X_reconstructed)
print("\nReconstruction error:")
print(torch.norm(X - X_reconstructed))
```

```
Original matrix:
tensor([[1., 2., 3.],
 [4., 5., 6.]])

Reconstructed matrix:
tensor([[1.0000, 2.0000, 3.0000],
 [4.0000, 5.0000, 6.0000]])

Reconstruction error:
tensor(1.8961e-06)
```

Figure: Demonstrate code for SVD.

## 4. SVD

### Definition

For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , there exists a factorization:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$$

where:

- $\mathbf{U} \in \mathbb{R}^{m \times m}$  is orthogonal ( $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_m$ )
- $\Sigma \in \mathbb{R}^{m \times n}$  is diagonal with non-negative entries
- $\mathbf{V} \in \mathbb{R}^{n \times n}$  is orthogonal ( $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_n$ )

## 4. SVD

### Definition

For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , there exists a factorization:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$$

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

$$\mathbf{U} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & \sqrt{5} & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{V}^\top = \begin{bmatrix} 0 & 0 & -1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3, \mathbf{u}_4], \Sigma = \text{diag}(3, \sqrt{5}, 2, 0), \mathbf{V}^\top = \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \\ \mathbf{v}_3^\top \\ \mathbf{v}_4^\top \\ \mathbf{v}_5^\top \end{bmatrix}$$

## 4. SVD

$$\mathbf{A} = \sum_{i=1}^4 \sigma_{ii} \mathbf{u}_i \mathbf{v}_i^\top$$

$$\mathbf{A} = 3 \begin{bmatrix} 0 \\ -1 \\ 0 \\ 0 \end{bmatrix} [0 \quad 0 \quad -1 \quad 0 \quad 0] + 3 \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix} [-\sqrt{0.2} \quad 0 \quad 0 \quad 0 \quad -\sqrt{0.8}] + 2 \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1 \end{bmatrix} [0 \quad -1 \quad 0 \quad 0 \quad 0]$$

## 4. SVD

- Compact SVD

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

$$\mathbf{U} = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -1 \end{bmatrix}, \Sigma = \begin{bmatrix} 3 & 0 & 0 \\ 0 & \sqrt{5} & 0 \\ 0 & 0 & 2 \end{bmatrix}, \mathbf{V}^\top = \begin{bmatrix} 0 & 0 & -1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \end{bmatrix}$$

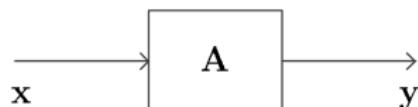
## 4. SVD

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

$$\mathbf{U} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & \sqrt{5} & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{V}^\top = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

$\Sigma$  (the singular values) is unique, while  $\mathbf{U}$  (left singular vectors) and  $\mathbf{V}$  (right singular vectors) are not unique.

## 4. SVD



(a) linear system



(b) equivalent system

**Figure:** Matrix-vector product from the perspective of SVD <sup>3</sup>.

<sup>3</sup><http://www.ee.cuhk.edu.hk/~wkma/>

# Power iteration

## Purpose

Power iteration is an algorithm to compute:

- The dominant eigenvalue  $\lambda_1$  (largest in magnitude)
- The corresponding eigenvector  $v_1$

## Key Assumption

- $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$
- Where  $\lambda_i$  are the eigenvalues of  $\mathbf{A}$
- Convergence rate depends on  $|\lambda_2/\lambda_1|$

## Algorithm description

**Algorithm 2** Power Iteration

- ```

1: Choose initial vector  $\mathbf{x}_0$  randomly, with  $\|\mathbf{x}_0\| = 1$ 
2: for  $k = 1, 2, \dots$  until convergence do
3:    $\mathbf{q}_k = \mathbf{A}\mathbf{x}_{k-1}$ 
4:    $\mathbf{x}_k = \mathbf{q}_k / \|\mathbf{q}_k\|$ 
5:    $\lambda_k = \mathbf{x}_k^\top \mathbf{A} \mathbf{x}_k$                                 ▷ Rayleigh quotient
6: end for

```

Power iteration

```
▶ #Demonstrate code for power interation
import torch

# Test matrix
X = torch.tensor([[4., 1.], [1., 3.]])

# Power iteration
v = torch.randn(2)
for _ in range(20):
    v = X @ v
    v = v / torch.norm(v)

eigenvalue = (v @ X @ v) / (v @ v)
print(f"\lambda={eigenvalue:.4f}, v={v}")
U,S,V = torch.svd(X)
print(S)
```

⇒ λ=4.6180, v=tensor([-0.8507, -0.5257])
tensor([4.6180, 2.3820])

Figure: Demonstrate code for power iteration.

4. SVD

Matrix-vector dot product

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

$$\mathbf{A} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0.3 & 0 & 0 & 0 \\ 0 & 0.2 & 0 & 0 \\ 0 & 0 & 0.1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 0 & -1 & 0 \\ -\sqrt{0.2} & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ -\sqrt{0.8} & 0 & 0 & 0 \end{bmatrix}$$

$$\mathbf{x} = [-\sqrt{0.2}, 0, 2, 0, -\sqrt{0.8}]^\top$$

$$\mathbf{y} = -2 \times 0.3 \times [0, -1, 0, 0]^\top + 1 \times 0.2 \times [-1, 0, 0, 0]^\top$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$
$$\mathbf{v}_1^\top \mathbf{x} \quad \sigma_1 \quad \mathbf{u}_1 \quad \mathbf{v}_2^\top \mathbf{x} \quad \sigma_2 \quad \mathbf{u}_2$$

4. SVD

Matrix-vector dot product

$$\mathbf{y} = \mathbf{Bx}$$

$$\mathbf{B} = \begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \end{bmatrix}, \mathbf{x} = [-\sqrt{0.2}, 0, 2, 0, -\sqrt{0.8}]^\top$$

$$\mathbf{U} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & \sqrt{5} & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{V}^\top = \begin{bmatrix} 0 & 0 & -1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

$$\mathbf{y} = -2 * 3 * [0, -1, 0, 0]^\top + 1 * \sqrt{5} * [-1, 0, 0, 0]^\top$$

4. SVD

Matrix-vector dot product

$$\mathbf{y} = \mathbf{Cx}$$

$$\mathbf{C} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 0 & -1 & 0 & 0 \\ -\sqrt{0.2} & 0 & 0 & 0 & -\sqrt{0.8} \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

$$\mathbf{y} = -2 * 3 * [0, -1, 0, 0]^T + 1 * 2 * [-1, 0, 0, 0]^T$$

Vector norms

$$\|\mathbf{x}\|_1 = \sum_i |x_i|$$

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^\top \mathbf{x}}$$

$$\|\mathbf{x}\|_p = \left[\sum_i |x_i|^p \right]^{1/p}$$

$$\|\mathbf{x}\|_\infty = \max_i |x_i|$$

Matrix norms

$$\|\mathbf{A}\|_1 = \max_j \sum_i |A_{ij}|$$

$$\|\mathbf{A}\|_2 = \sqrt{\max \text{eig}(\mathbf{A}^\top \mathbf{A})}$$

$$\|\mathbf{A}\|_p = \left(\max_{\|\mathbf{x}\|_p=1} \|\mathbf{Ax}\|_p \right)$$

$$\|\mathbf{A}\|_\infty = \max_i \sum_j |A_{ij}|$$

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2} = \sqrt{\text{trace}(\mathbf{A}^\top \mathbf{A})}$$

Inequalities

Assuming \mathbf{A} is an $m \times n$, and $d = \text{rank}(\mathbf{A})$

	$\ \mathbf{A}\ _{\max}$	$\ \mathbf{A}\ _1$	$\ \mathbf{A}\ _\infty$	$\ \mathbf{A}\ _2$	$\ \mathbf{A}\ _F$	$\ \mathbf{A}\ _{\text{KF}}$
$\ \mathbf{A}\ _{\max}$	1	1	1	1	1	1
$\ \mathbf{A}\ _1$	m	1	m	\sqrt{m}	\sqrt{m}	\sqrt{m}
$\ \mathbf{A}\ _\infty$	n	n	1	\sqrt{n}	\sqrt{n}	\sqrt{n}
$\ \mathbf{A}\ _2$	\sqrt{mn}	\sqrt{n}	\sqrt{m}	1	1	1
$\ \mathbf{A}\ _F$	\sqrt{mn}	\sqrt{n}	\sqrt{m}	\sqrt{d}	1	1
$\ \mathbf{A}\ _{\text{KF}}$	\sqrt{mnd}	\sqrt{nd}	\sqrt{md}	d	\sqrt{d}	1

which are to be read as, e.g.

$$\|\mathbf{A}\|_2 \leq \sqrt{m} \cdot \|\mathbf{A}\|_\infty$$

Theorem: Singular Value Bounds of a Gaussian Random Matrix

Theorem

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ have i.i.d. standard Gaussian entries, i.e., $\mathbf{A}_{i,j} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$. For every $m \geq n$, we have the following inequality:

$$\sqrt{m} - \sqrt{n} \leq \mathbb{E}[\sigma_{\min}(\mathbf{A})] \leq \mathbb{E}[\sigma_{\max}(\mathbf{A})] \leq \sqrt{m} + \sqrt{n},$$

where $\sigma_{\min}(\mathbf{A})$ and $\sigma_{\max}(\mathbf{A})$ denote the minimal and maximal singular values, respectively.

Lemma: Upper Bound of Weight Matrix Norm (Xavier Initialization)

Lemma

Let $\mathbf{A} \in \mathbb{R}^{n_{in} \times n_{out}}$ have i.i.d. standard Gaussian entries, i.e.,

$$A_{i,j} \stackrel{iid}{\sim} \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right).$$

We have the following inequality for its maximum singular value:

$$\mathbb{E}[\sigma_{\max}(\mathbf{A})] \leq 2.$$

$$A_{i,j} \stackrel{iid}{\sim} \mathcal{N}\left(0, \left(\frac{1}{\sqrt{n_{in}} + \sqrt{n_{out}}}\right)^2\right)$$

Expected norm of addition of two unit random vectors

$\mathbf{x}_1 = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$, \mathbf{x} is a random high-dimensional vector.

$\mathbf{y}_1 = \frac{\mathbf{y}}{\|\mathbf{y}\|_2}$, \mathbf{y} is a random high-dimensional vector.

$$\mathbf{z}_1 = \mathbf{x}_1 + \mathbf{y}_1$$

$$\mathbb{E}(\|\mathbf{z}\|_2) = \sqrt{2}$$

Expected norm of addition of two unit random vectors

```
▶ #Demonstrate code for power interation
import torch
# Set parameters
dim = 1024      # dimension of vectors
n_trials = 1000   # number of trials
# Store norms of results
result_norms = torch.zeros(n_trials)
# Perform trials
for i in range(n_trials):
    # Create two random vectors
    v1 = torch.randn(dim)
    v2 = torch.randn(dim)

    # L2 normalize both vectors
    v1_normalized = v1 / torch.norm(v1)
    v2_normalized = v2 / torch.norm(v2)

    # Add normalized vectors and compute norm
    result = v1_normalized + v2_normalized
    result_norms[i] = torch.norm(result)

# Calculate statistics
mean_norm = torch.mean(result_norms)
var_norm = torch.var(result_norms)

print(f"Statistics for {n_trials} trials with {dim}-dimensional vectors:")
print(f"Mean norm: {mean_norm:.6f}")
print(f"Variance of norm: {var_norm:.6f}")
print(f"Standard deviation of norm: {torch.sqrt(var_norm):.6f}")

➡ Statistics for 1000 trials with 1024-dimensional vectors:
Mean norm: 1.413116
Variance of norm: 0.000480
Standard deviation of norm: 0.021916
```

Figure: Demonstrate code for expected norm of addition of two unit vectors.

Expected norm under orthogonal transformation

- Given: $\mathbf{y} = \mathbf{Ax}$ where \mathbf{A} is orthogonal
- Properties of orthogonal matrix \mathbf{A} :

$$\mathbf{A}^\top \mathbf{A} = \mathbf{AA}^\top = \mathbf{I}$$

- Consider the norm of \mathbf{y} :

$$\begin{aligned}\|\mathbf{y}\|^2 &= \mathbf{y}^\top \mathbf{y} = (\mathbf{Ax})^\top (\mathbf{Ax}) \\ &= \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|^2\end{aligned}$$

- Taking expectations:

$$\mathbb{E}[\|\mathbf{y}\|^2] = \mathbb{E}[\|\mathbf{x}\|^2]$$

- Therefore:

$$\mathbb{E}[\|\mathbf{y}\|] = \mathbb{E}[\|\mathbf{x}\|]$$

Expected Norm with Orthogonal Transformation

```
▶ import torch

# Parameters
n = 100 # dimension
trials = 1000
norms = torch.zeros(trials)
# Create orthogonal matrix using QR decomposition
A = torch.randn(n, n)
Q, _ = torch.linalg.qr(A) # Q is orthogonal
for i in range(trials):
    # Random input vector
    x = torch.randn(n)
    x = x / torch.norm(x)

    # Compute norm of y = Ax
    norms[i] = torch.norm(Q @ x)

print(f"Expected norm: {torch.mean(norms):.4f} ± {torch.std(norms):.4f}")
print(f"Input vector norm (mean): {torch.mean(torch.tensor(norms[i])):.4f}")
```

→ Expected norm: 1.0000 ± 0.0000
Input vector norm (mean): 1.0000

Figure: Expected Norm with L2 Normalized Columns.

Expected norm with L2 normalized columns

- Given: $\mathbf{y} = \mathbf{Ax}$ where $\mathbf{A} \in \mathbb{R}^{n \times n}$ with L2 normalized columns
- Properties of \mathbf{A} :

$$\|\mathbf{a}_i\|_2 = 1 \text{ for all columns } i$$

$$\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$$

- Consider the norm squared of \mathbf{y} :

$$\begin{aligned}\|\mathbf{y}\|^2 &= \mathbf{y}^\top \mathbf{y} = (\mathbf{Ax})^\top (\mathbf{Ax}) = \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} \\ &= \mathbf{x}^\top \mathbf{Mx} \text{ where } \mathbf{M} = \mathbf{A}^\top \mathbf{A}\end{aligned}$$

- Properties of $\mathbf{M} = \mathbf{A}^\top \mathbf{A}$:
 - $M_{ii} = 1$ (due to L2 normalization)
 - $M_{ij} = \mathbf{a}_i^\top \mathbf{a}_j$ (inner products of columns)
- Therefore:

$$\begin{aligned}\mathbb{E}[\|\mathbf{y}\|^2] &= \mathbb{E}[\mathbf{x}^\top \mathbf{Mx}] \\ &= \mathbb{E}[\|\mathbf{x}\|^2] + \sum_{i \neq j} M_{ij} \mathbb{E}[x_i x_j]\end{aligned}$$

Expected Norm with L2 Normalized Columns

```
▶ import torch

# Parameters
n = 100 # dimension of square matrix
trials = 1000
norms = torch.zeros(trials)
# Create matrix A with L2 normalized columns
A = torch.randn(n, n)
A = A / torch.norm(A, dim=0, keepdim=True) # normalize columns
for i in range(trials):
    # Random input vector
    x = torch.randn(n)
    x = x / torch.norm(x)

    # Compute norm of y = Ax
    norms[i] = torch.norm(A @ x)

print(f"Expected norm: {torch.mean(norms):.4f} ± {torch.std(norms):.4f}")
U,S,V = torch.svd(A)
print('Expected norm: ', torch.sqrt(torch.pow(S, 2).mean()))
```

→ Expected norm: 0.9944 ± 0.0689
Expected norm: tensor(1.)

Figure: Expected Norm with L2 Normalized Columns.

Expected norm with L2 normalized rows

- Given: $\mathbf{y} = \mathbf{A}\mathbf{x}$ where:
- $\mathbf{A} \in \mathbb{R}^{m \times n}$ with L2 normalized rows
- Properties of \mathbf{A} : $\|\mathbf{a}_i^\top\|_2 = 1$ for all rows i

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \mathbf{a}_2^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}$$

- Consider $\|\mathbf{y}\|^2$: $\|\mathbf{y}\|^2 = \sum_{i=1}^m y_i^2 = \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 = \sum_{i=1}^m \mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}$
- $\mathbb{E}[\|\mathbf{y}\|^2] = \sum_{i=1}^m \mathbb{E}[\mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}] = \sum_{i=1}^m (\mathbb{E}[\|\mathbf{x}\|^2] \cdot \frac{1}{n} + \text{correlation terms})$

Note:

- Each row contributes $\|\mathbf{x}\|^2/n$ in expectation if \mathbf{x} is isotropic
- Total contribution depends on row correlations

Expected Norm with L2 Normalized Rows

```
▶ #Demonstrate code to Expected Norm with L2 Normalized Rows (Rectangle)
import torch

# Parameters
m, n = 100, 100
trials = 1000
norms = torch.zeros(trials)
A = torch.randn(m, n)
A = A / torch.norm(A, dim=1, keepdim=True)
for i in range(trials):
    # Create and normalize A and x
    x = torch.randn(n)
    x = x / torch.norm(x)

    # Compute norm of y = Ax
    norms[i] = torch.norm(A @ x)

print(f"Expected norm: {torch.mean(norms):.4f} ± {torch.std(norms):.4f}")
U,S,V = torch.svd(A)
print('Expected norm: ', torch.sqrt(torch.pow(S, 2).mean()))
```

→ Expected norm: 0.9981 ± 0.0699
Expected norm: tensor(1.0000)

Figure: Expected Norm with L2 Normalized Rows (Rectangle).

Expected norm with L2 normalized row (underdetermined)

- Given: $\mathbf{y} = \mathbf{A}\mathbf{x}$ where:
- $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n$ and L2 normalized rows
- Properties of \mathbf{A} : $\|\mathbf{a}_i^\top\|_2 = 1$ for all rows i

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \mathbf{a}_2^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}$$

- Consider $\|\mathbf{y}\|^2$: $\|\mathbf{y}\|^2 = \sum_{i=1}^m y_i^2 = \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 = \sum_{i=1}^m \mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}$
- $\mathbb{E}[\|\mathbf{y}\|^2] = \sum_{i=1}^m \mathbb{E}[\mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}] = \sum_{i=1}^m (\mathbb{E}[\|\mathbf{x}\|^2] \cdot \frac{1}{n} + \text{correlation terms})$
Note:
 - Each row contributes $\|\mathbf{x}\|^2/n$ in expectation if \mathbf{x} is isotropic
 - Total contribution depends on row correlations

Expected Norm with L2 Normalized Rows

```
▶ #Demonstrate code to Expected Norm with L2 Normalized Rows (Underdetermined)
import torch

# Parameters
m, n = 20, 100
trials = 1000
norms = torch.zeros(trials)
A = torch.randn(m, n)
A = A / torch.norm(A, dim=1, keepdim=True)
for i in range(trials):
    # Create and normalize A and x
    x = torch.randn(n)
    x = x / torch.norm(x)

    # Compute norm of y = Ax
    norms[i] = torch.norm(A @ x)

print(f"Expected norm: {torch.mean(norms):.4f} ± {torch.std(norms):.4f}")
U,S,V = torch.svd(A)
print('Expected norm: ', torch.sqrt(torch.pow(S, 2).mean()))
```

→ Expected norm: 0.4389 ± 0.0697
Expected norm: tensor(1.)

Figure: Expected Norm with L2 Normalized Rows (Underdetermined).

Expected norm with L2 normalized row (overdetermined)

- Given: $\mathbf{y} = \mathbf{Ax}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m > n$ and L2 normalized rows
Thus, we have

$$\|\mathbf{a}_i^\top\|_2 = 1 \text{ for all rows } i, \mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \mathbf{a}_2^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}$$

- Consider $\|\mathbf{y}\|_2^2$: $\|\mathbf{y}\|_2^2 = \sum_{i=1}^m y_i^2 = \sum_{i=1}^m (\mathbf{a}_i^\top \mathbf{x})^2 = \sum_{i=1}^m \mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}$
- Taking expectations:

$$\mathbb{E}[\|\mathbf{y}\|^2] = \sum_{i=1}^m \mathbb{E}[\mathbf{x}^\top (\mathbf{a}_i \mathbf{a}_i^\top) \mathbf{x}] = \sum_{i=1}^m \left(\mathbb{E}[\|\mathbf{x}\|^2] \cdot \frac{1}{n} + \text{correlation terms} \right)$$

- Key observations:
 - $\mathbb{E}[\|\mathbf{y}\|^2]$ typically larger than $\mathbb{E}[\|\mathbf{x}\|^2]$ due to redundancy

Expected Norm with L2 Normalized Rows

```
▶ #Demonstrate code to Expected Norm with L2 Normalized Rows
import torch

# Parameters
m, n = 100, 20
trials = 1000
norms = torch.zeros(trials)
A = torch.randn(m, n)
A = A / torch.norm(A, dim=1, keepdim=True)
for i in range(trials):
    # Create and normalize A and x
    x = torch.randn(n)
    x = x / torch.norm(x)

    # Compute norm of y = Ax
    norms[i] = torch.norm(A @ x)

print(f"Expected norm: {torch.mean(norms):.4f} ± {torch.std(norms):.4f}")
U,S,V = torch.svd(A)
print('Expected norm: ', torch.sqrt(torch.pow(S, 2).mean()))
```

→ Expected norm: 2.2343 ± 0.1432
Expected norm: tensor(2.2361)

Figure: Expected Norm with L2 Normalized Rows (Overdetermined).

Eigenvalue decomposition & SVD: a close relationship

- **EVD vs SVD:**

EVD: $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$, where \mathbf{A} is $n \times n$

SVD: $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$, where \mathbf{A} is $m \times n$

- **Key Connections:**

- ▶ For symmetric matrices \mathbf{A} :

$$\mathbf{A}\mathbf{A}^\top = (\mathbf{U}\Sigma\mathbf{V}^\top)(\mathbf{V}\Sigma\mathbf{U}^\top) = \mathbf{U}\Sigma^2\mathbf{U}^\top$$

- ▶ Singular values (σ_i) relate to eigenvalues (λ_i):

$$\sigma_i = \sqrt{\lambda_i(\mathbf{A}\mathbf{A}^\top)}$$

- **Important Implications:**

- ▶ Left singular vectors (\mathbf{U}) are eigenvectors of $\mathbf{A}\mathbf{A}^\top$
 - ▶ Right singular vectors (\mathbf{V}) are eigenvectors of $\mathbf{A}^\top\mathbf{A}$

Part 2: Matrix Calculus

[7, 9, 16, 10, 11]

Single-variable calculus⁴

Basic Properties and Formulas

If $f(x)$ and $g(x)$ are differentiable functions (the derivative exists), c and n are any real numbers,

$$1. \frac{d}{dx}(c) = 0$$

$$4. (f(x) \pm g(x))' = f'(x) \pm g'(x)$$

$$2. (c f(x))' = c f'(x)$$

$$5. (f(x) g(x))' = f'(x) g(x) + f(x) g'(x) \text{ - Product Rule}$$

$$3. \frac{d}{dx}(x^n) = n x^{n-1} \text{ - Power Rule}$$

$$6. \left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x) g(x) - f(x) g'(x)}{(g(x))^2} \text{ - Quotient Rule}$$

$$7. \frac{d}{dx}\left(f(g(x))\right) = f'(g(x)) g'(x) \text{ - Chain Rule}$$

⁴<https://tutorial.math.lamar.edu>

Single-variable calculus

Common Derivatives

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(\sin(x)) = \cos(x)$$

$$\frac{d}{dx}(\cos(x)) = -\sin(x)$$

$$\frac{d}{dx}(\tan(x)) = \sec^2(x)$$

$$\frac{d}{dx}(\sec(x)) = \sec(x)\tan(x)$$

$$\frac{d}{dx}(\csc(x)) = -\csc(x)\cot(x)$$

$$\frac{d}{dx}(\cot(x)) = -\csc^2(x)$$

$$\frac{d}{dx}(\sin^{-1}(x)) = \frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\cos^{-1}(x)) = -\frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\tan^{-1}(x)) = \frac{1}{1+x^2}$$

$$\frac{d}{dx}(a^x) = a^x \ln(a)$$

$$\frac{d}{dx}(\mathbf{e}^x) = \mathbf{e}^x$$

$$\frac{d}{dx}(\ln(x)) = \frac{1}{x}, \quad x > 0$$

$$\frac{d}{dx}(\ln|x|) = \frac{1}{x}, \quad x \neq 0$$

$$\frac{d}{dx}(\log_a(x)) = \frac{1}{x \ln(a)}, \quad x > 0$$

Single-variable calculus

Chain Rule Variants

The chain rule applied to some specific functions.

$$1. \frac{d}{dx} \left([f(x)]^n \right) = n [f(x)]^{n-1} f'(x)$$

$$2. \frac{d}{dx} \left(e^{f(x)} \right) = f'(x) e^{f(x)}$$

$$3. \frac{d}{dx} \left(\ln [f(x)] \right) = \frac{f'(x)}{f(x)}$$

$$4. \frac{d}{dx} \left(\sin [f(x)] \right) = f'(x) \cos [f(x)]$$

$$5. \frac{d}{dx} \left(\cos [f(x)] \right) = -f'(x) \sin [f(x)]$$

$$6. \frac{d}{dx} \left(\tan [f(x)] \right) = f'(x) \sec^2 [f(x)]$$

$$7. \frac{d}{dx} \left(\sec [f(x)] \right) = f'(x) \sec [f(x)] \tan [f(x)]$$

$$8. \frac{d}{dx} \left(\tan^{-1} [f(x)] \right) = \frac{f'(x)}{1 + [f(x)]^2}$$

Single-variable calculus

Higher Order Derivatives

The 2nd Derivative is denoted as

$$f''(x) = f^{(2)}(x) = \frac{d^2f}{dx^2} \text{ and is defined as}$$

$f''(x) = (f'(x))'$, i.e. the derivative of the first derivative, $f'(x)$.

The n^{th} Derivative is denoted as

$$f^{(n)}(x) = \frac{d^n f}{dx^n} \text{ and is defined as}$$

$f^{(n)}(x) = (f^{(n-1)}(x))'$, i.e. the derivative of the $(n-1)^{st}$ derivative, $f^{(n-1)}(x)$.

- Consider the chain rule expression:

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} \stackrel{?}{=} \frac{dy}{dx} \frac{dz}{dy}$$

- Two main conventions:

① Numerator Layout

- Elements of derivative arranged by numerator
- $[\frac{dz}{dx}]_{ij} = \frac{\partial z_i}{\partial x_j}$
- Chain rule: left-to-right multiplication

② Denominator Layout

- Elements of derivative arranged by denominator
- $[\frac{dz}{dx}]_{ij} = \frac{\partial z_j}{\partial x_i}$
- Chain rule: right-to-left multiplication

- Choice affects matrix multiplication order in chain rule!

⁵ Matrix calculus wiki

Numerator layout in matrix calculus

- Case 1: Scalar by Vector (y scalar, \mathbf{x} vector)

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}$$

- Case 2: Vector by Scalar (\mathbf{y} vector, x scalar)

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$$

Numerator layout in matrix calculus

- Case 3: Vector by Vector (\mathbf{y} vector, \mathbf{x} vector)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

- Case 4: Scalar by Matrix (y scalar, \mathbf{X} matrix)

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \cdots & \frac{\partial y}{\partial x_{p1}} \\ \frac{\partial y}{\partial x_{12}} & \frac{\partial y}{\partial x_{22}} & \cdots & \frac{\partial y}{\partial x_{p2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1q}} & \frac{\partial y}{\partial x_{2q}} & \cdots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

Denominator layout in matrix calculus

- Case 1: Scalar by Vector (y scalar, \mathbf{x} vector)

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

- Case 2: Vector by Scalar (\mathbf{y} vector, x scalar)

$$\frac{\partial \mathbf{y}}{\partial x} = \left[\frac{\partial y_1}{\partial x} \quad \frac{\partial y_2}{\partial x} \quad \cdots \quad \frac{\partial y_m}{\partial x} \right]$$

Denominator layout in matrix calculus

- Case 3: Vector by Vector (\mathbf{y} vector, \mathbf{x} vector)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \frac{\partial y_2}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

- Case 4: Scalar by Matrix (y scalar, \mathbf{X} matrix)

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{12}} & \dots & \frac{\partial y}{\partial x_{1q}} \\ \frac{\partial y}{\partial x_{21}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{2q}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{p1}} & \frac{\partial y}{\partial x_{p2}} & \dots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

Derivatives and Gradients

- For a scalar-valued function $f(\cdot) : \mathbf{x} \rightarrow f(\mathbf{x})$

- The derivative is a $1 \times m$ row vector:

$$\frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_m} \right]$$

- The gradient $\nabla_{\mathbf{x}} f$ is the transposed column vector:

$$\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}} \right)^{\top} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_m} \end{bmatrix}$$

Jacobian matrix

- For a vector-valued function $\mathbf{f}(\cdot) : \mathbf{x} \rightarrow \mathbf{f}(\mathbf{x})$
- The derivative, known as the Jacobian, is an $n \times m$ matrix
- The (i,j) -th element is the scalar derivative of the i -th output w.r.t the j -th input:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_m} \end{bmatrix}$$

Scalar by matrix derivatives

- For a scalar-valued function $f(\cdot) : \begin{matrix} \mathbf{X} \\ m \times n \end{matrix} \rightarrow f(\mathbf{X}) \begin{matrix} \\ 1 \times 1 \end{matrix}$
- The derivative is an $m \times n$ matrix:

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial X_{1,1}} & \cdots & \frac{\partial f}{\partial X_{m,1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial X_{1,n}} & \cdots & \frac{\partial f}{\partial X_{m,n}} \end{bmatrix}$$

- The gradient has dimensions matching the input matrix: $m \times n$

Matrix-by-scalar derivatives

- For a matrix-valued function $\text{vec}(\cdot) : \underset{1 \times 1}{x} \rightarrow \underset{p \times q}{\mathbf{F}(x)}$
- The derivative is a $p \times q$ matrix:

$$\frac{\partial \mathbf{F}}{\partial x} = \begin{bmatrix} \frac{\partial F_{1,1}}{\partial x} & \dots & \frac{\partial F_{1,q}}{\partial x} \\ \vdots & \ddots & \vdots \\ \frac{\partial F_{p,1}}{\partial x} & \dots & \frac{\partial F_{p,q}}{\partial x} \end{bmatrix}$$

- Each element is the scalar derivative of the corresponding matrix element

Vector-by-matrix derivatives

- For a vector-valued function $\text{vec}(\cdot) : \underset{m \times n}{\mathbf{X}} \rightarrow \underset{p \times 1}{\mathbf{f}(\mathbf{X})}$
- The derivative is a 3rd-order tensor with dimensions $p \times m \times n$:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial X_{1,1}} & \cdots & \frac{\partial \mathbf{f}}{\partial X_{m,1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{f}}{\partial X_{1,n}} & \cdots & \frac{\partial \mathbf{f}}{\partial X_{m,n}} \end{bmatrix}$$

- Each element $\frac{\partial \mathbf{f}}{\partial X_{i,j}}$ is a p -dimensional vector
- This represents a stack of p matrices, each of size $m \times n$

Matrix-by-vector derivatives

- For a matrix-valued function $\text{vec}(\cdot) : \mathbf{x}_{m \times 1} \rightarrow \mathbf{F}(\mathbf{x})_{p \times q}$
- The derivative is a 3rd-order tensor with dimensions $p \times q \times m$:

$$\frac{\partial \mathbf{F}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{F}}{\partial x_1} & \frac{\partial \mathbf{F}}{\partial x_2} & \cdots & \frac{\partial \mathbf{F}}{\partial x_m} \end{bmatrix}$$

- Each element $\frac{\partial \mathbf{F}}{\partial x_i}$ is a $p \times q$ matrix
- The tensor can be viewed as m matrices of size $p \times q$

Tensor derivatives

- For a vector-valued function $\mathbf{f}(\cdot) : \mathbb{X}_{m \times n} \rightarrow \mathbb{F}(\mathbb{X})_{p \times q}$
- The derivative is a 4th-order tensor with dimensions $p \times q \times m \times n$

$$\frac{\partial \mathbf{F}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial \mathbf{F}}{\partial X_{1,1}} & \cdots & \frac{\partial \mathbf{F}}{\partial X_{m,1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{F}}{\partial X_{1,n}} & \cdots & \frac{\partial \mathbf{F}}{\partial X_{m,n}} \end{bmatrix}$$

- This represents $p \times q$ matrices of size $m \times n$, one for each output dimension

Gradient and Hessian⁶

The oracle is

$$f = \mathbf{x}^\top \mathbf{A}\mathbf{x} + \mathbf{b}^\top \mathbf{x}$$

The gradient (first-order derivative) is

$$\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}} \right)^\top = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} + \mathbf{b}$$

The Hessian (second-order derivative) is

$$\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \mathbf{A} + \mathbf{A}^\top$$

⁶Matrix calculus calculator

Gradient and Hessian

The oracle is

$$\begin{aligned} f &= (\mathbf{A}\mathbf{x} + \mathbf{b})^\top(\mathbf{C}\mathbf{x} + \mathbf{d}) \\ &= \mathbf{x}^\top \mathbf{A}^\top \mathbf{C}\mathbf{x} + \mathbf{x}^\top \mathbf{A}^\top \mathbf{d} + \mathbf{b}^\top \mathbf{C}\mathbf{x} + \mathbf{b}^\top \mathbf{d} \end{aligned}$$

The gradient (first-order derivative) is

$$\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}} \right)^\top = (\mathbf{A}^\top \mathbf{C} + \mathbf{C}^\top \mathbf{A})\mathbf{x} + \mathbf{A}^\top \mathbf{d} + \mathbf{C}^\top \mathbf{b}$$

The Hessian (second-order derivative) is

$$\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \mathbf{A}^\top \mathbf{C} + \mathbf{C}^\top \mathbf{A}$$

Gradient example: L2 norm

- Problem: Find $\nabla_{\mathbf{x}} f$ where $f(\mathbf{x}) = \|\mathbf{x} - \mathbf{a}\|_2$
- Note: $\|\mathbf{x} - \mathbf{a}\|_2 = \sqrt{(\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}$ is scalar
- Using chain rule:

$$\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial \sqrt{(\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}}{\partial (\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})} \cdot \frac{\partial (\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}{\partial \mathbf{x}}$$

- Where:
 - ▶ First term = $\frac{1}{2\sqrt{(\mathbf{x}-\mathbf{a})^\top (\mathbf{x}-\mathbf{a})}}$
 - ▶ Second term = $2(\mathbf{x} - \mathbf{a})^\top$

- Final result:

$$\frac{\partial f}{\partial \mathbf{x}} = \frac{(\mathbf{x} - \mathbf{a})^\top}{\sqrt{(\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}}$$

Derivative of normalized vector

- Problem: Find $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ where $\mathbf{y} = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$
- Using quotient rule and chain rule:

$$\begin{aligned}\frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \frac{\partial}{\partial \mathbf{x}} \left(\frac{\mathbf{x}}{\sqrt{\mathbf{x}^\top \mathbf{x}}} \right) \\ &= \frac{\mathbf{I} \sqrt{\mathbf{x}^\top \mathbf{x}} - \mathbf{x} \frac{\partial}{\partial \mathbf{x}} \sqrt{\mathbf{x}^\top \mathbf{x}}}{(\sqrt{\mathbf{x}^\top \mathbf{x}})^2}\end{aligned}$$

- Note that $\frac{\partial}{\partial \mathbf{x}} \sqrt{\mathbf{x}^\top \mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$

- Final result:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{1}{\|\mathbf{x}\|_2} \left(\mathbf{I} - \frac{\mathbf{x} \mathbf{x}^\top}{\mathbf{x}^\top \mathbf{x}} \right)$$

- This is a projection matrix scaled by $\frac{1}{\|\mathbf{x}\|_2}$

Norm derivative with respect to x

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. At the top, there is an input field for the derivative expression: $(x-a)/\text{norm2}(x-a)$. To its right is a dropdown menu for "w.r.t." containing "x" and a dropdown arrow. Below this, the derivative is calculated as:

$$\frac{\partial}{\partial x} \left((x - a) / \|x - a\|_2 \right) = 1/t_1 \cdot \mathbb{I} - 1/t_1^3 \cdot t_0 \cdot t_0^\top$$

Underneath the derivative, the text "where" is followed by two definitions:

$$t_0 = x - a$$
$$t_1 = \|t_0\|_2$$

Further down, another "where" section is shown with two dropdown menus:

- "a is a" dropdown set to "vector"
- "x is a" dropdown set to "vector"

On the right side of the interface, there are several buttons and settings:

- "Export functions as" button with options "Python" (selected) and "Latex".
- "Common subexpressions" button with the value "ON".

Figure: Expected Norm with L2 Normalized Columns.

Derivative of least-square

- Problem: Find $\frac{\partial f}{\partial \mathbf{x}}$ where $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|^2$
- Expanding the squared norm:

$$\begin{aligned}f(\mathbf{x}) &= (\mathbf{Ax} - \mathbf{y})^\top (\mathbf{Ax} - \mathbf{y}) \\&= \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} - \mathbf{x}^\top \mathbf{A}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{Ax} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

- Taking the derivative:

$$\frac{\partial f}{\partial \mathbf{x}} = 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{A} - 2\mathbf{y}^\top \mathbf{A}$$

- Final result can also be written as:

$$\frac{\partial f}{\partial \mathbf{x}} = 2(\mathbf{Ax} - \mathbf{y})^\top \mathbf{A}$$

- This is the gradient used in linear least squares optimization

Least-square derivative with respect to y

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

derivative of w.r.t.

$$\frac{\partial}{\partial y} ((A \cdot x - y)^\top \cdot (A \cdot x - y)) = -2 \cdot (A \cdot x - y)$$

where

A is a

x is a

y is a

Export functions as

Common subexpressions

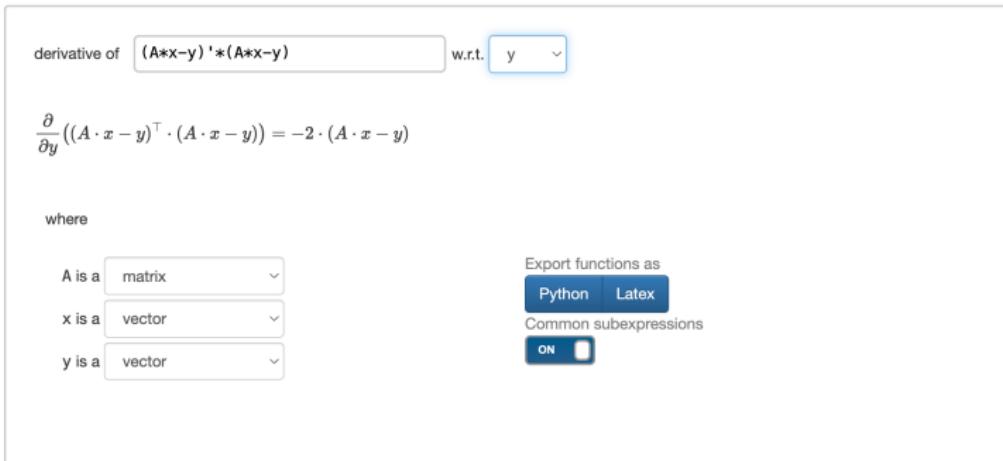


Figure: Expected Norm with L2 Normalized Columns.

Least-square derivative with respect to x

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. In the top input field, the expression $(A \cdot x - y) \cdot (A \cdot x - y)$ is entered. Below it, the text "w.r.t." is followed by a dropdown menu with options "x" and "y", where "x" is selected. The resulting derivative is displayed as $\frac{\partial}{\partial x} ((A \cdot x - y)^\top \cdot (A \cdot x - y)) = 2 \cdot A^\top \cdot (A \cdot x - y)$. To the left of the derivative, the text "where" is followed by three dropdown menus: "A is a matrix", "x is a vector", and "y is a vector". On the right side, there are buttons for "Export functions as Python" and "Latex", and a switch for "Common subexpressions" which is set to "ON".

Figure: Expected Norm with L2 Normalized Columns.

Least-square derivative with respect to A

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

derivative of w.r.t.

$$\frac{\partial}{\partial A} \left((A \cdot x - y)^T \cdot (A \cdot x - y) \right) = 2 \cdot (A \cdot x - y) \cdot x^T$$

where

A is a

x is a

y is a

Export functions as

Common subexpressions
 ON

Figure: Expected Norm with L2 Normalized Columns.

Extended examples of matrix calculus

- Basic vector operations:

- $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = 2\mathbf{x}$
- $f(\mathbf{x}) = \|\mathbf{x}\|_2$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$
- $f(\mathbf{x}) = \|\mathbf{x}\|_2^2$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = 2\mathbf{x}$

- Vector transformations:

- $\mathbf{y} = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}$, $\left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right) = \frac{1}{\|\mathbf{x}\|_2} (\mathbf{I} - \frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^\top \mathbf{x}})$
- $f(\mathbf{x}) = \|\mathbf{x} - \mathbf{a}\|_2$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = \frac{(\mathbf{x} - \mathbf{a})}{\|\mathbf{x} - \mathbf{a}\|_2}$

- Matrix-vector operations:

- $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A}\mathbf{x}$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = (\mathbf{A} + \mathbf{A}^\top)\mathbf{x}$
- $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = 2\mathbf{A}^\top(\mathbf{A}\mathbf{x} - \mathbf{y})$
- $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x}\|_2$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = \frac{\mathbf{A}^\top \mathbf{A}\mathbf{x}}{\|\mathbf{A}\mathbf{x}\|_2}$

- Common objective functions:

- $f(\mathbf{x}) = \log(1 + e^{\mathbf{x}^\top \mathbf{a}})$, $\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top = \frac{e^{\mathbf{x}^\top \mathbf{a}}}{1+e^{\mathbf{x}^\top \mathbf{a}}} \mathbf{a}$
- $f(\mathbf{x}) = -\sum_i y_i \log(x_i)$, $\frac{\partial f}{\partial \mathbf{x}} = -[\frac{y_1}{x_1}, \dots, \frac{y_n}{x_n}]^\top$

Vec operator

- For matrix product \mathbf{AB} :
 $\mathbf{A} \in \mathcal{R}^{m \times n}, \mathbf{B} \in \mathcal{R}^{n \times k}$

$$\text{vec}(\mathbf{AB}) = (\mathbf{I}_k \otimes \mathbf{A}) \text{vec}(\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{I}_m) \text{vec}(\mathbf{A})$$

- Properties:
 - \otimes denotes Kronecker product
 - \mathbf{I}_k is $k \times k$ identity matrix
 - \mathbf{I}_m is $m \times m$ identity matrix

Vec operator derivatives: two matrices

- Let $\mathbf{P} = \mathbf{AB}$ where:

- $\mathbf{A} \in \mathcal{R}^{m \times n}$
- $\mathbf{B} \in \mathcal{R}^{n \times k}$

- Derivatives:

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{A})} = \mathbf{B}^\top \otimes \mathbf{I}_m$$

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{B})} = \mathbf{I}_k \otimes \mathbf{A}$$

Vec operator: triple matrices

- Let $\mathbf{P} = \mathbf{ABC}$ where:

- $\mathbf{A} \in \mathcal{R}^{m \times n}$
- $\mathbf{B} \in \mathcal{R}^{n \times k}$
- $\mathbf{C} \in \mathcal{R}^{k \times l}$

$$\begin{aligned}\text{vec}(\mathbf{ABC}) &= (\mathbf{C}^\top \mathbf{B}^\top \otimes \mathbf{I}_m) \text{vec}(\mathbf{A}) \\ &= (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B}) \\ &= (\mathbf{I}_l \otimes \mathbf{AB}) \text{vec}(\mathbf{C})\end{aligned}$$

Vec operator derivatives: triple matrices

- Let $\mathbf{P} = \mathbf{ABC}$ where:

- $\mathbf{A} \in \mathcal{R}^{m \times n}$
- $\mathbf{B} \in \mathcal{R}^{n \times k}$
- $\mathbf{C} \in \mathcal{R}^{k \times l}$

- Derivatives:

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{A})} = \mathbf{C}^\top \mathbf{B}^\top \otimes \mathbf{I}_m$$

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{B})} = \mathbf{C}^\top \otimes \mathbf{A}$$

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{C})} = \mathbf{I}_l \otimes \mathbf{AB}$$

First-order matrix derivatives: basic cases

- Basic trace derivatives:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}) = \mathbf{I}$$

- With constant matrix \mathbf{A} :

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}\mathbf{A}) = \mathbf{A}^\top$$

- With transpose variations:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^\top \mathbf{A}) = \mathbf{A}$$

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{A}\mathbf{X}^\top) = \mathbf{A}$$

First-order matrix derivatives: advanced cases

- With two constant matrices \mathbf{A}, \mathbf{B} :

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{AXB}) = \mathbf{A}^T \mathbf{B}^T$$

- Transpose variations:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{AX}^T \mathbf{B}) = \mathbf{BA}$$

- Kronecker product:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{A} \otimes \mathbf{X}) = \text{trace}(\mathbf{A}) \mathbf{I}$$

\mathbf{AXB} derivative respect to \mathbf{X}

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. In the top input field, the expression $\text{tr}(\mathbf{A} \cdot \mathbf{X} \cdot \mathbf{B})$ is entered. To its right, the text "w.r.t." is followed by a dropdown menu containing the letter "X". Below this, the result of the differentiation is displayed as:

$$\frac{\partial}{\partial \mathbf{X}} (\text{tr}(\mathbf{A} \cdot \mathbf{X} \cdot \mathbf{B})) = (\mathbf{B} \cdot \mathbf{A})^\top$$

Underneath the result, the word "where" is followed by three dropdown menus: "A is a matrix", "B is a matrix", and "X is a matrix", each with a "matrix" option selected. To the right of these dropdowns are two buttons: "Export functions as Python" (which is highlighted in blue) and "Latex". Further down, there is a button labeled "Common subexpressions" with a switch labeled "ON".

Figure: $\text{trace}(\mathbf{AXB})$ derivative respect to \mathbf{X} .

Second-order matrix derivatives: basic cases

- Quadratic term:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^2) = 2\mathbf{X}^\top$$

- With single constant matrix \mathbf{B} :

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^2 \mathbf{B}) = (\mathbf{X}\mathbf{B} + \mathbf{B}\mathbf{X})^\top$$

- Equivalent forms:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^\top \mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}\mathbf{X}^\top) = 2\mathbf{X}$$

- With quadratic term:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^\top \mathbf{B}\mathbf{X}) = \mathbf{B}\mathbf{X} + \mathbf{B}^\top \mathbf{X}$$

Second-order matrix derivatives: advanced cases

- Various forms with $\mathbf{X}\mathbf{X}^\top$:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{B}\mathbf{X}\mathbf{X}^\top) = \mathbf{B}\mathbf{X} + \mathbf{B}^\top \mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}\mathbf{X}^\top \mathbf{B}) = \mathbf{B}\mathbf{X} + \mathbf{B}^\top \mathbf{X}$$

- Forms with $\mathbf{X}\mathbf{B}\mathbf{X}^\top$:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}\mathbf{B}\mathbf{X}^\top) = \mathbf{X}\mathbf{B}^\top + \mathbf{X}\mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{B}\mathbf{X}^\top \mathbf{X}) = \mathbf{X}\mathbf{B}^\top + \mathbf{X}\mathbf{B}$$

- Double matrix product:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}) = \mathbf{A}^\top \mathbf{X}^\top \mathbf{B}^\top + \mathbf{B}^\top \mathbf{X}^\top \mathbf{A}^\top$$

trace(BXX^T) derivative respect to X

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. In the top input field, the expression $\text{tr}(B \cdot X \cdot X^T)$ is entered. Below it, the variable X is selected from a dropdown menu labeled "w.r.t.". The resulting derivative is displayed as $\frac{\partial}{\partial X} (\text{tr}(B \cdot X \cdot X^T)) = B^T \cdot X + B \cdot X$. To the left, under "where", there are two dropdown menus: one for B set to "matrix" and one for X also set to "matrix". On the right, there are buttons for "Export functions as" (Python, Latex), a "Common subexpressions" toggle switch (set to ON), and a small search icon.

Figure: $\text{trace}(BXX^T)$ derivative respect to X .

Higher-order matrix derivatives: k-th power cases

- General k-th power rule:

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^k) = k(\mathbf{X}^{k-1})^\top$$

- With constant matrix \mathbf{A} :

$$\frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{AX}^k) = \sum_{r=0}^{k-1} (\mathbf{X}^r \mathbf{AX}^{k-r-1})^\top$$

- Special cases:

- $k = 2: \frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^2) = 2\mathbf{X}^\top$
- $k = 3: \frac{\partial}{\partial \mathbf{X}} \text{trace}(\mathbf{X}^3) = 3(\mathbf{X}^2)^\top$

Higher-order matrix derivatives: complex cases

- Complex trace derivative with multiple matrices:

$$\begin{aligned}\frac{\partial}{\partial \mathbf{X}} \text{trace}[\mathbf{B}^T \mathbf{X}^T \mathbf{C} \mathbf{X} \mathbf{X}^T \mathbf{C} \mathbf{X} \mathbf{B}] = \\ \mathbf{C} \mathbf{X} \mathbf{X}^T \mathbf{C} \mathbf{X} \mathbf{B} \mathbf{B}^T \\ + \mathbf{C}^T \mathbf{X} \mathbf{B} \mathbf{B}^T \mathbf{X}^T \mathbf{C}^T \mathbf{X} \\ + \mathbf{C} \mathbf{X} \mathbf{B} \mathbf{B}^T \mathbf{X}^T \mathbf{C} \mathbf{X} \\ + \mathbf{C}^T \mathbf{X} \mathbf{X}^T \mathbf{C}^T \mathbf{X} \mathbf{B} \mathbf{B}^T\end{aligned}$$

Theorem: differential product rule

Theorem (Differential Product Rule)

Let \mathbf{A}, \mathbf{B} be two matrices. Then, we have the differential product rule for \mathbf{AB} :

$$d(\mathbf{AB}) = (d\mathbf{A})\mathbf{B} + \mathbf{A}(d\mathbf{B})$$

- By the differential of matrix \mathbf{A} , we think of it as a small (unconstrained) change in matrix \mathbf{A}
- Later, constraints may be placed on the allowed perturbations

Differential operator: scalar case

- For a scalar function $f(x)$, the differential is:

$$df = f(x + dx) - f(x)$$

- This represents an infinitesimal change in f due to an infinitesimal change in x
- Key properties:
 - $d(f + g) = df + dg$
 - $d(fg) = (df)g + f(dg)$ (product rule)
 - $d(f/g) = \frac{(df)g - f(dg)}{g^2}$ (quotient rule)

Differential operator: vector case

- For a vector \mathbf{x} , the differential is:

$$d\mathbf{x} = \begin{bmatrix} dx_1 \\ dx_2 \\ \vdots \\ dx_n \end{bmatrix}$$

- Properties extend component-wise:
 - $d(\mathbf{x} + \mathbf{y}) = d\mathbf{x} + d\mathbf{y}$
 - $d(\mathbf{x}^\top \mathbf{y}) = (d\mathbf{x})^\top \mathbf{y} + \mathbf{x}^\top (d\mathbf{y})$
 - $d\|\mathbf{x}\| = \frac{\mathbf{x}^\top d\mathbf{x}}{\|\mathbf{x}\|}$
- The differential represents a small displacement in vector space

Differential operator: matrix case

- For a matrix \mathbf{A} , the differential is:

$$d\mathbf{A} = \begin{bmatrix} dA_{11} & dA_{12} & \cdots & dA_{1n} \\ dA_{21} & dA_{22} & \cdots & dA_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ dA_{m1} & dA_{m2} & \cdots & dA_{mn} \end{bmatrix}$$

- Key matrix differential rules:

- $d(\mathbf{AB}) = (d\mathbf{A})\mathbf{B} + \mathbf{A}(d\mathbf{B})$
- $d(\mathbf{A}^\top) = (d\mathbf{A})^\top$
- $d(\mathbf{A}^{-1}) = -\mathbf{A}^{-1}(d\mathbf{A})\mathbf{A}^{-1}$

Differential operator: tensor case

- For a tensor \mathbf{A} , the differential applies element-wise:

$$dA_{ijkl\dots} = A_{ijkl\dots}(x + dx) - A_{ijkl\dots}(x)$$

- Properties:

- ▶ Linearity: $d(\alpha\mathbf{A} + \beta\mathbf{B}) = \alpha d\mathbf{A} + \beta d\mathbf{B}$
- ▶ Product Rule extends to tensor products
- ▶ Maintains tensor order and structure

- Applications:

- ▶ Stress/strain analysis in continuum mechanics
- ▶ Multi-dimensional data analysis
- ▶ Higher-order derivatives in deep learning

Differential operator: basic rules

- Constant rule:

$$\partial \mathbf{A} = 0$$

- Linear operations:

$$\partial(\alpha \mathbf{X}) = \alpha \partial \mathbf{X}$$

$$\partial(\mathbf{X} + \mathbf{Y}) = \partial \mathbf{X} + \partial \mathbf{Y}$$

- Product rules:

$$\partial(\mathbf{X} \mathbf{Y}) = (\partial \mathbf{X}) \mathbf{Y} + \mathbf{X} (\partial \mathbf{Y})$$

$$\partial(\mathbf{X} \circ \mathbf{Y}) = (\partial \mathbf{X}) \circ \mathbf{Y} + \mathbf{X} \circ (\partial \mathbf{Y})$$

$$\partial(\mathbf{X} \otimes \mathbf{Y}) = (\partial \mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes (\partial \mathbf{Y})$$

- Transpose and Hermitian:

$$\partial \mathbf{X}^\top = (\partial \mathbf{X})^\top$$

$$\partial \mathbf{X}^H = (\partial \mathbf{X})^H$$

Differential operator: advanced rules

- Trace operation:

$$\partial(\text{trace}(\mathbf{X})) = \text{trace}(\partial\mathbf{X})$$

- Matrix inverse:

$$\partial(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}(\partial\mathbf{X})\mathbf{X}^{-1}$$

- Determinant (two forms):

$$\partial(\det(\mathbf{X})) = \text{trace}(\text{adj}(\mathbf{X})\partial\mathbf{X})$$

$$\partial(\det(\mathbf{X})) = \det(\mathbf{X}) \text{trace}(\mathbf{X}^{-1}\partial\mathbf{X})$$

- Log-determinant:

$$\partial(\ln(\det(\mathbf{X}))) = \text{trace}(\mathbf{X}^{-1}\partial\mathbf{X})$$

- Note: These rules are essential in:

- ▶ Matrix optimization
- ▶ Neural network backpropagation
- ▶ Statistical parameter estimation

Matrix trace differentiates

$$\begin{aligned} d \operatorname{trace}(\mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C}) &= d \operatorname{trace}(\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T) = \operatorname{trace}(d(\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T)) \\ &= \operatorname{trace}(\mathbf{C} \mathbf{A} \mathbf{X} d(\mathbf{B} \mathbf{X}^T) + d(\mathbf{C} \mathbf{A} \mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}(\mathbf{C} \mathbf{A} \mathbf{X} d(\mathbf{B} \mathbf{X}^T)) + \operatorname{trace}(d(\mathbf{C} \mathbf{A} \mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}(\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} d(\mathbf{X}^T)) + \operatorname{trace}(\mathbf{C} \mathbf{A} (d\mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}(\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} (d\mathbf{X})^T) + \operatorname{trace}(\mathbf{C} \mathbf{A} (d\mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}((\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} (d\mathbf{X})^T)^T) + \operatorname{trace}(\mathbf{C} \mathbf{A} (d\mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}((d\mathbf{X}) \mathbf{B}^T \mathbf{X}^T \mathbf{A}^T \mathbf{C}^T) + \operatorname{trace}(\mathbf{C} \mathbf{A} (d\mathbf{X}) \mathbf{B} \mathbf{X}^T) \\ &= \operatorname{trace}(\mathbf{B}^T \mathbf{X}^T \mathbf{A}^T \mathbf{C}^T (d\mathbf{X})) + \operatorname{trace}(\mathbf{B} \mathbf{X}^T \mathbf{C} \mathbf{A} (d\mathbf{X})) \\ &= \operatorname{trace}((\mathbf{B}^T \mathbf{X}^T \mathbf{A}^T \mathbf{C}^T + \mathbf{B} \mathbf{X}^T \mathbf{C} \mathbf{A}) d\mathbf{X}) \\ &= \operatorname{trace}((\mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B} + \mathbf{A}^T \mathbf{C}^T \mathbf{X} \mathbf{B}^T)^T d\mathbf{X}) \end{aligned}$$

Lipschitz continuity

Definition

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz continuous if there exists $L \geq 0$ such that:

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. The smallest such L is called the Lipschitz constant.

Linear transformations

For a linear transformation $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$\|\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{y}\| = \|\mathbf{A}(\mathbf{x} - \mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$$

Matrix norms and Lipschitz constants

Key Relationship

For a matrix \mathbf{A} , its Lipschitz constant equals its operator norm:

$$L = \|\mathbf{A}\| = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|}$$

Matrix norms

- Spectral norm (2-norm): $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}$
- Frobenius norm: $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |a_{ij}|^2}$
- Maximum absolute row sum: $\|\mathbf{A}\|_\infty = \max_i \sum_j |a_{ij}|$
- Maximum absolute column sum: $\|\mathbf{A}\|_1 = \max_j \sum_i |a_{ij}|$

Properties

Basic Properties

For matrices \mathbf{A} and \mathbf{B} :

- $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$
- $\|c\mathbf{A}\| = |c| \|\mathbf{A}\|$ for scalar c
- $\|\mathbf{A}\| \geq 0$ with equality iff $\mathbf{A} = 0$
- $\|\mathbf{A}\| \geq \rho(\mathbf{A})$ (spectral radius)

Relation to singular values

For the spectral norm:

$$\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$$

where σ_1 is the largest singular value of \mathbf{A}

Applications in optimization

Gradient Descent

For function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{A}^\top \mathbf{A}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$:

- Gradient: $\nabla f(\mathbf{x}) = \mathbf{A}^\top \mathbf{A}\mathbf{x} - \mathbf{b}$
- Lipschitz constant: $L = \|\mathbf{A}^\top \mathbf{A}\|_2 = \sigma_1^2(\mathbf{A})$
- Step size bound: $\alpha < \frac{2}{L}$ for convergence

Convergence Analysis

Rate of convergence depends on condition number:

$$\kappa(\mathbf{A}) = \frac{\sigma_1(\mathbf{A})}{\sigma_n(\mathbf{A})} = \frac{L}{\mu}$$

where μ is strong convexity parameter

Practical computation

Computing Lipschitz Constants

For matrix \mathbf{A} :

- ① SVD method: $L = \sigma_1(\mathbf{A})$
- ② Power iteration for $\mathbf{A}^\top \mathbf{A}$
- ③ Direct norm computation for simple cases

Example

For 2×2 matrix:

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

Spectral norm (Lipschitz constant):

$$L = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} = \sqrt{\frac{\alpha + \sqrt{\alpha^2 - 4\beta}}{2}}$$

where $\alpha = a^2 + b^2 + c^2 + d^2$ and $\beta = (ad - bc)^2$

Part 3: Transformer

[25, 17, 18, 12, 19, 20, 29, 27, 2]

Transformer

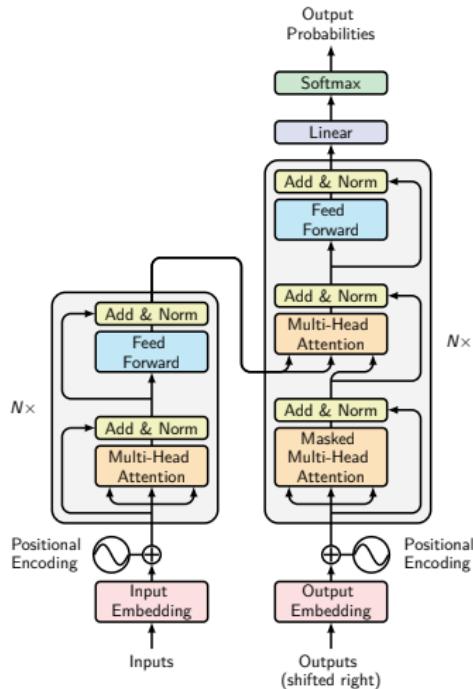
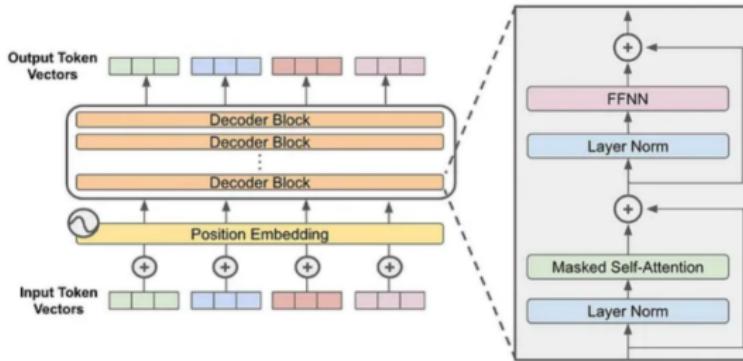


Figure: Original Transformer

Pure-decoder



$$\text{Softmax} \left[\frac{\text{Query Matrix} \times \text{Transposed Key Matrix}}{\sqrt{d}} \right] = \begin{array}{c} \text{Unnormalized Attention Scores} \\ \begin{array}{|c|c|c|c|c|} \hline 1.1 & 2.3 & 0.5 & 2.1 & -1.2 \\ \hline -0.1 & 0.5 & 1.3 & -0.1 & 0.1 \\ \hline 2.3 & 0.2 & 3.3 & -1.0 & -0.4 \\ \hline 0.3 & 1.2 & -0.3 & 0.0 & 1.4 \\ \hline 0.1 & -2.9 & -1.1 & -4.2 & 0.4 \\ \hline \end{array} \\ \xrightarrow{\quad\quad\quad} \text{Masked Attention Scores} \\ \begin{array}{|c|c|c|c|c|} \hline 1.1 & \cancel{-0.1} & \cancel{0.5} & \cancel{2.1} & \cancel{-1.2} \\ \hline -0.1 & 0.5 & \cancel{1.3} & \cancel{-0.1} & 0.1 \\ \hline 2.3 & 0.2 & 3.3 & \cancel{-1.0} & \cancel{-0.4} \\ \hline 0.3 & 1.2 & -0.3 & 0.0 & 1.4 \\ \hline 0.1 & -2.9 & -1.1 & -4.2 & 0.4 \\ \hline \end{array} \end{array}$$

Figure: GPT.

Encoder

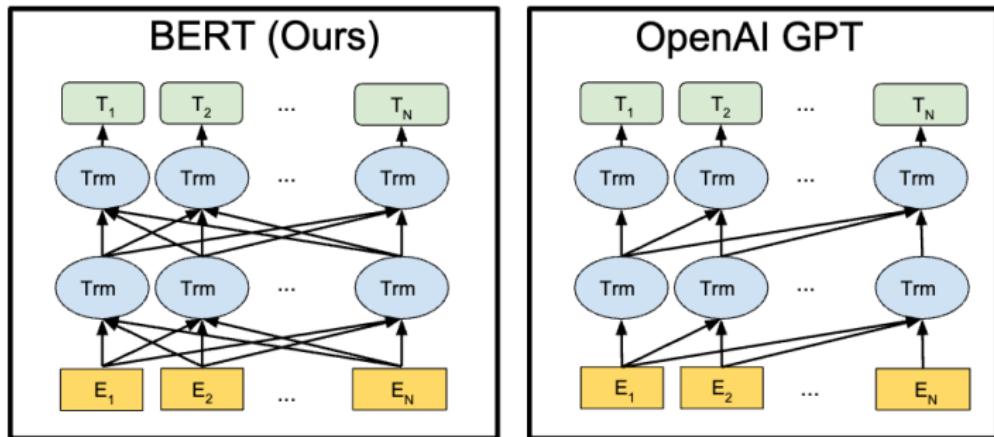


Figure: BERT vs GPT.

Pure-encoder

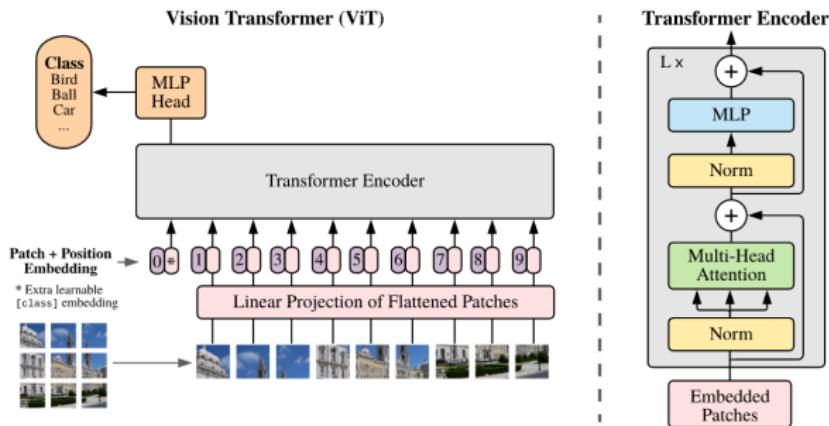
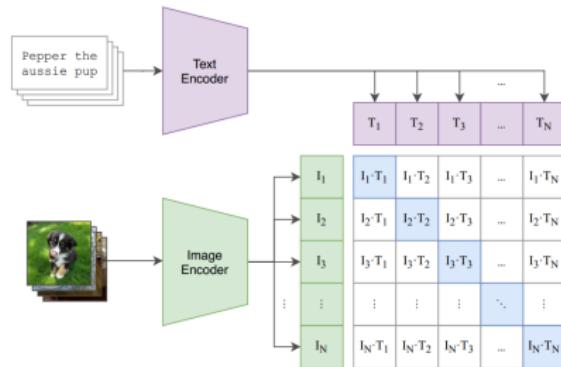


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by [Vaswani et al. \(2017\)](#).

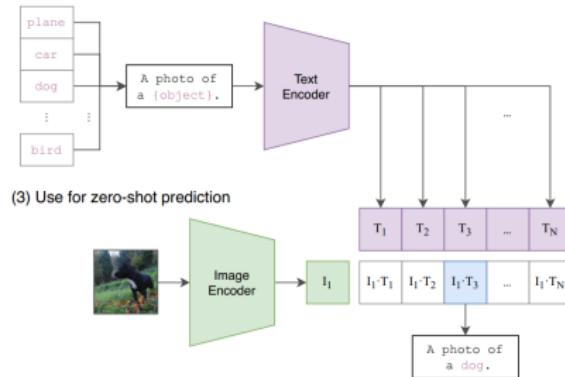
Figure: Vision Transformer.

Pure-encoder

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

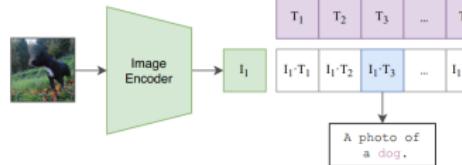


Figure: CLIP.

Encoder-decoder

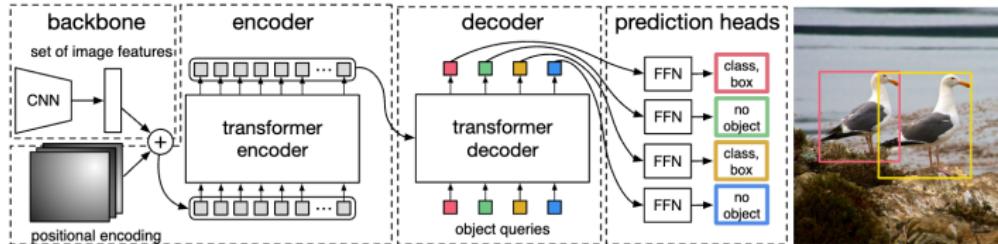


Fig. 2: DETR uses a conventional CNN backbone to learn a 2D representation of an input image. The model flattens it and supplements it with a positional encoding before passing it into a transformer encoder. A transformer decoder then takes as input a small fixed number of learned positional embeddings, which we call *object queries*, and additionally attends to the encoder output. We pass each output embedding of the decoder to a shared feed forward network (FFN) that predicts either a detection (class and bounding box) or a “no object” class.

Figure: DETR.

Number of parameters

If L is the number of layer, d is the hidden size, V is the number of vocabulary.

- $L \times (d \times d * 12 + (4 + 5 + 2 * 2) * d) + V * d$

For example, nanoGPT-small:

- $12 * (768 * 768 * 12 + 13 * 768) + 50000 * 768 = 123454464$

Memory usage in training

Let $\mathbf{X} \in \mathcal{R}^{b \times n \times d}$, L is the number of layer, h is the number of heads.

In the forward process, we need to cache the following variables

- $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{O} \in \mathcal{R}^{b \times n \times d}$
- $\mathbf{P}, \mathbf{A} \in \mathcal{R}^{b \times n \times n}$
- in each block, we have h times of \mathbf{P} and \mathbf{A}
- in FFN, we have $\mathcal{R}^{b \times n \times 4d}$ (ReLU replace can be used.)
- in residual shortcut, we have two $\mathcal{R}^{b \times n \times d}$

Note:

Flash-attention is extremely important in reducing memory cost and speeding up the training.

Computation complexity

- Self-attention module

$$\mathcal{O}(L \times b \times n \times d \times 4d + L \times b \times n \times n \times d \times 2)$$

- FFN module

$$\mathcal{O}(L \times b \times n \times d \times 8d)$$

Positional encoding ⁷

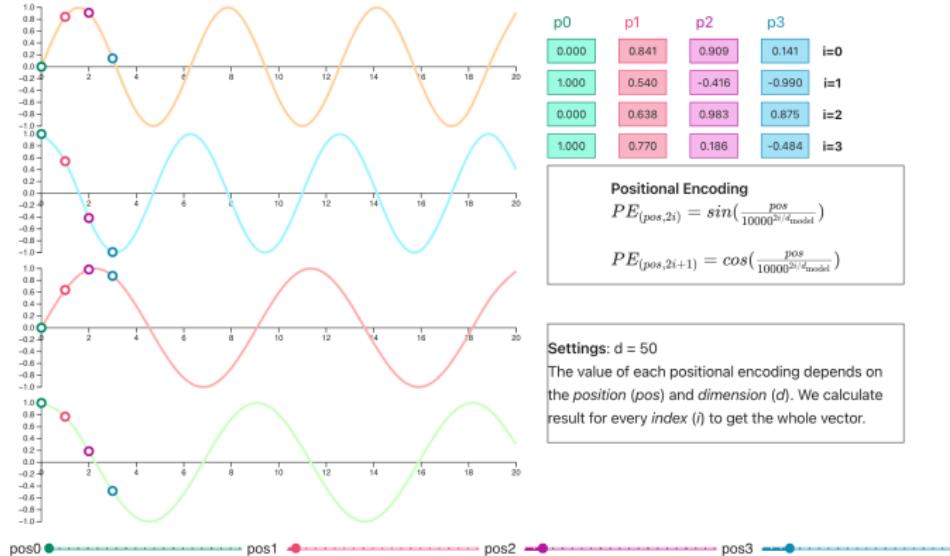
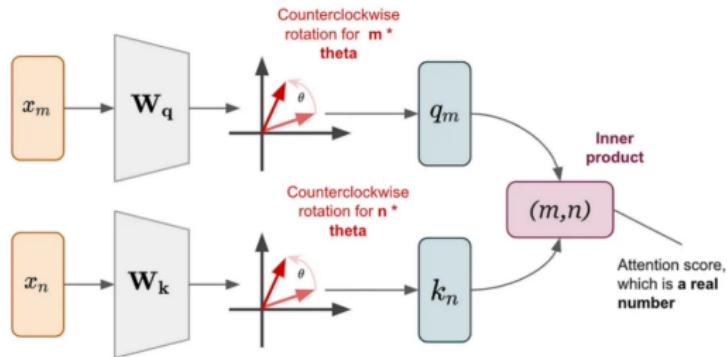


Figure: Visualization of position embedding.

Rotary position embedding⁸



$$f_q(\mathbf{x}_m, m) = (\mathbf{W}_q \mathbf{x}_m) e^{im\theta}$$

$$f_k(\mathbf{x}_n, n) = (\mathbf{W}_k \mathbf{x}_n) e^{in\theta}$$

$$g(\mathbf{x}_m, \mathbf{x}_n, m - n) = \text{Re}[(\mathbf{W}_q \mathbf{x}_m)(\mathbf{W}_k \mathbf{x}_n)^* e^{i(m-n)\theta}]$$

Figure: Visualization of RoPE.

⁸RoPE blog

Rotary position embedding

$$\mathbf{R}_{\Theta,m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$

And the overall transformation applied to a token embedding is:

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta,m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m$$

$$\mathbf{R}_{\Theta,m}^d \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_d \\ x_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$

Figure: Computational process of RoPE.

Advanced skills

- Flash-attention
- GQA
- Multi-head Latent Attention
- KV-cache
- MoE
- LoRA

Chain Rule

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} &= \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \\&= \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} \cdots \frac{\partial \mathbf{x}^L}{\partial \mathbf{x}^{L-1}} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^L} \\&= \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} \cdots \frac{\partial \mathbf{x}^L}{\partial \mathbf{x}^{L-1}} \frac{\partial \mathbf{o}}{\partial \mathbf{x}^L} \frac{\partial \mathcal{L}}{\partial \mathbf{o}} \\&= \left(\prod_{k=l+1}^L \frac{\partial \mathbf{x}^k}{\partial \mathbf{x}^{k-1}} \right) \frac{\partial \mathbf{o}}{\partial \mathbf{x}^L} \frac{\partial \mathcal{L}}{\partial \mathbf{o}} \\&= \left(\prod_{k=l+1}^L \begin{pmatrix} \frac{\partial \mathbf{h}^k}{\partial \mathbf{x}^{k-1}} & \frac{\partial \mathbf{r}^k}{\partial \mathbf{h}^k} & \frac{\partial \mathbf{x}^k}{\partial \mathbf{r}^k} \end{pmatrix} \right) \frac{\partial \mathbf{o}}{\partial \mathbf{x}^L} \frac{\partial \mathcal{L}}{\partial \mathbf{o}} \\&= \left(\prod_{k=l+1}^L \left(\mathbf{W}^k \right)^\top \text{diag} \left(\mathbf{1} \left(\mathbf{W}^k \mathbf{x}^{k-1} > 0 \right) \right) \frac{\sqrt{D} \left(\mathbf{I} - \frac{1}{D} \mathbf{1} \mathbf{1}^\top \right) \left(\mathbf{I} - \frac{\mathbf{y}^k \mathbf{y}^{k\top}}{\|\mathbf{y}^k\|_2^2 + \epsilon} \right) \text{diag} \left(\boldsymbol{\gamma}^k \right)}{\sqrt{\|\mathbf{y}^k\|_2^2 + \epsilon}} \right) \mathbf{O}^\top (\mathbf{p} - \mathbf{t})\end{aligned}$$

Forward process of Transformer:

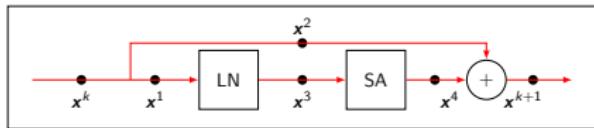


Figure: Forward process of self-attention module

$$x^1 = x^2 = x^k,$$

$$x^3 = \text{LN}(x^1),$$

$$x^4 = \text{SA}(x^3),$$

$$x^{k+1} = x^2 + x^4.$$

Backward process of Transformer:

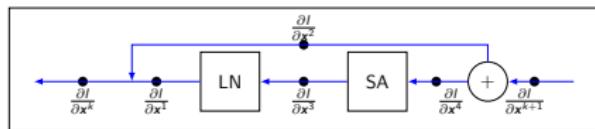


Figure: Backward process of self-attention module

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{x}^{k+1}} &= \frac{\partial l}{\partial \mathbf{x}^4} = \frac{\partial l}{\partial \mathbf{x}^2}, \\ \frac{\partial l}{\partial \mathbf{x}^3} &= \frac{\partial l}{\partial \mathbf{x}^4} J_{x^4}(x^3), \\ \frac{\partial l}{\partial \mathbf{x}^1} &= \frac{\partial l}{\partial \mathbf{x}^4} J_{x^4}(x^3), \\ \frac{\partial l}{\partial \mathbf{x}^1} &= \frac{\partial l}{\partial \mathbf{x}^3} J_{x^3}(x^1), \\ \frac{\partial l}{\partial \mathbf{x}^k} &= \frac{\partial l}{\partial \mathbf{x}^2} + \frac{\partial l}{\partial \mathbf{x}^1},\end{aligned}$$

$J_{x^3}(x^1)$ is the Jacobian matrix of x^3 with respect to x^1 .

Denote it as $J_{x^3}(x^1) = \frac{\partial x^3}{\partial x^1}$.

LayerNorm

- Forward Process

$$\text{LN}(\mathbf{x}) = \gamma \odot \frac{\sqrt{d}\mathbf{y}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}} + \beta, \quad \text{and } \mathbf{y} = \left(\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top \right) \mathbf{x},$$

where $\mathbf{x} \in \mathcal{R}^d$

- Backward Process

$$\frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}} \text{Diag}(\gamma) \left(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^\top}{\|\mathbf{y}\|_2^2 + \epsilon} \right) \left(\mathbf{I} - \frac{1}{d}\mathbf{1}\mathbf{1}^\top \right).$$

$$\frac{\partial \text{LN}(\mathbf{x})}{\partial \gamma} = \text{Diag}(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^\top}{\|\mathbf{y}\|_2^2 + \epsilon}), \quad \frac{\partial \text{LN}(\mathbf{x})}{\partial \beta} = \frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{x}},$$

$$\frac{\partial \text{LN}(\mathbf{x})}{\partial \mathbf{y}} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{y}\|_2^2 + \epsilon}} \text{Diag}(\gamma) \left(\mathbf{I} - \frac{\mathbf{y}\mathbf{y}^\top}{\|\mathbf{y}\|_2^2 + \epsilon} \right)$$

Feed-forward network

- ▷ Forward Process

$$\mathbf{y} = \text{FFN}(\mathbf{x}) = \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x}),$$

- ▷ Backward Process

$$\mathbf{J}_y(\mathbf{x}) = \frac{\partial \text{FFN}(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{W}_2 \text{Diag}(\mathbf{1}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 > 0)) \mathbf{W}_1$$

Residual block

- ▷ Forward Process

$$\mathbf{y} = \mathbf{x} + f(\mathbf{x})$$

- ▷ Backward Process

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{I} + \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

Softmax and cross-entropy loss

▷ Forward Process

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

$$\mathcal{L} = - \sum_{i=1}^n y_i \log(p_i)$$

where \mathbf{y} is one-hot label vector

▷ Backward Process

$$\nabla_{\mathbf{x}} \mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \right)^{\top} = \mathbf{p} - \mathbf{y} = \text{softmax}(\mathbf{x}) - \mathbf{y}$$

- Note: The gradient simplifies elegantly to $\mathbf{p} - \mathbf{y}$ due to the properties of softmax and cross-entropy.

Softmax and cross-entropy loss

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. At the top, there is a text input field containing the expression `-sum(y.*log(exp(x)/sum(exp(x))))` and a dropdown menu set to "w.r.t. x". Below this, the derivative is calculated:

$$\frac{\partial}{\partial x} (-\text{sum}(y \odot \log(\exp(x) / \text{sum}(\exp(x))))) = \\ -(1/t_1 \cdot y \odot t_0 \oslash t_2 - 1/t_1^2 \cdot t_0^\top \cdot (y \oslash t_2) \cdot t_0)$$

Below the derivative, the text "where" is followed by three definitions:

$$t_0 = \exp(x)$$
$$t_1 = \text{sum}(t_0)$$
$$t_2 = 1/\text{sum}(t_0) \cdot t_0$$

Further down, another "where" section is shown with two dropdown menus: "x is a" set to "vector" and "y is a" also set to "vector". To the right of these are buttons for "Export functions as Python" and "Latex", and a switch for "Common subexpressions" which is set to "ON".

Figure: Gradient of softmax and cross-entropy.

Self-attention: two forms

- Form 1: if $\mathbf{X} \in \mathcal{R}^{d \times n}$,

$$\mathbf{Y} = SA(\mathbf{X}) = \mathbf{W}_v \mathbf{X} \text{softmax}\left(\frac{\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}}{\sqrt{d/h}}\right)$$

where softmax is conducted in the row direction, $\mathbf{Y} \in \mathcal{R}^{d \times n}$.

- Form 2: if $\mathbf{X} \in \mathcal{R}^{n \times d}$,

$$\mathbf{Y} = SA(\mathbf{X}) = \text{softmax}\left(\frac{\mathbf{X} \mathbf{W}_q \mathbf{W}_k^\top \mathbf{X}^\top}{\sqrt{d/h}}\right) \mathbf{X} \mathbf{W}_v$$

where softmax is conducted in the column direction, $\mathbf{Y} \in \mathcal{R}^{n \times d}$.

Kronecker product of self-attention

Let $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}$, $\mathbf{A} = \text{softmax}\left(\frac{\mathbf{P}}{\sqrt{d_q}}\right)$,

where \mathbf{A} is the attention map, and $\frac{\mathbf{P}}{\sqrt{d_q}}$ is usually called as the logit.

Assume $\mathbf{X} \in \mathcal{R}^{d \times n}$, $\mathbf{W}_q \in \mathcal{R}^{d_q \times d}$, $\mathbf{W}_k \in \mathcal{R}^{d_q \times d}$, according to vectorization and matrix calculus, we have the following derivations:

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top \mathbf{W}_k)} = \mathbf{X}^\top \otimes \mathbf{X}^\top, \quad \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{X}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{K} + (\mathbf{I}_n \otimes \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k),$$

$$\frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_q^\top)} = (\mathbf{W}_k \mathbf{X})^\top \otimes \mathbf{X}^\top, \quad \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{W}_k)} = \mathbf{X}^\top \otimes (\mathbf{W}_q \mathbf{X})^\top,$$

where $\mathbf{I}_n \in \mathcal{R}^{n \times n}$ denotes an identity matrix with shape $n \times n$, \mathbf{K} is the commutation matrix, which depends on the dimensions of \mathbf{X} . Since $\mathbf{X} \in \mathcal{R}^{d \times n}$, then we know $\mathbf{K} \in \mathcal{R}^{nd \times nd}$. The commutation matrix \mathbf{K} has the property that $\text{vec}(\mathbf{X}^\top) = \mathbf{K} \text{vec}(\mathbf{X})$ for any matrix \mathbf{X} .

Single-head self-attention

▷ Forward Process

$$\mathbf{Y} = \mathbf{W}_v \mathbf{X} \mathbf{A},$$

where $\mathbf{P} = \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{X}$, $\mathbf{A} = \text{softmax}\left(\frac{\mathbf{P}}{\sqrt{d_q}}\right)$.

▷ Backward Process

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{X})}, \\ &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})}.\end{aligned}$$

Kronecker product of softmax operator

Define

$$\mathbf{B} = \exp\left(\frac{\mathbf{P}}{\sqrt{d_q}}\right), \quad \mathbf{C} = \mathbf{1}\mathbf{1}^\top \exp\left(\frac{\mathbf{P}}{\sqrt{d_q}}\right) = \mathbf{1}\mathbf{1}^\top \mathbf{B}.$$

In this way, we can represent the attention matrix \mathbf{A} as $\mathbf{A} = \mathbf{B} \oslash \mathbf{C}$.
Then, we can vectorize the equation,

$$\text{vec}(\mathbf{A}) = \text{vec}(\mathbf{B} \oslash \mathbf{C}) = \text{vec}(\mathbf{B}) \oslash \text{vec}(\mathbf{C}).$$

According to the chain rule, we have

$$\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} = \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{B})} \frac{\partial \text{vec}(\mathbf{B})}{\partial \text{vec}(\mathbf{P})} + \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{C})} \frac{\partial \text{vec}(\mathbf{C})}{\partial \text{vec}(\mathbf{P})}.$$

Kronecker product of softmax operator

Let us calculate each individual term. We have

$$\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{B})} = \mathbf{1}_{n^2} \oslash \text{Diag}(\text{vec}(\mathbf{C})),$$

$$\frac{\partial \text{vec}(\mathbf{B})}{\partial \text{vec}(\mathbf{P})} = \frac{\text{Diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}},$$

$$\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{C})} = -\text{Diag}(\text{vec}(\mathbf{B}) \oslash (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))),$$

$$\frac{\partial \text{vec}(\mathbf{C})}{\partial \text{vec}(\mathbf{P})} = \frac{(\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{Diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}},$$

where \mathbf{I}_n is a $n \times n$ identity matrix and \otimes is the Kronecker product, $\mathbf{1}_{nn}$ denotes a \mathcal{R}^{n^2} all one vector.

Kronecker product of softmax operator

Substitute these four term into the chain rule, we have

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} &= \frac{\left(\mathbf{1}_{n^2} \oslash \text{diag}(\text{vec}(\mathbf{C}))\right) \text{diag}(\text{vec}(\mathbf{B})) - \text{diag}(\text{vec}(\mathbf{B}) \oslash (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))) (\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}} \\ &= \frac{\text{diag}(\text{vec}(\mathbf{A})) - \text{diag}(\text{vec}(\mathbf{B}) \oslash (\text{vec}(\mathbf{C}) \odot \text{vec}(\mathbf{C}))) (\mathbf{I}_n \otimes \mathbf{1}_n \mathbf{1}_n^\top) \text{diag}(\text{vec}(\mathbf{B}))}{\sqrt{d_q}} \\ &= \frac{\text{blockdiag}(\text{diag}(\mathbf{A}_{:,1}) - \mathbf{A}_{:,1} \mathbf{A}_{:,1}^\top, \dots, \text{diag}(\mathbf{A}_{:,n}) - \mathbf{A}_{:,n} \mathbf{A}_{:,n}^\top)}{\sqrt{d_q}}.\end{aligned}$$

If \mathbf{A} and \mathbf{P} are two vectors, if we use \mathbf{a} and \mathbf{p} denote them individually, then we know

$$\frac{\partial \text{vec}(\mathbf{a})}{\partial \text{vec}(\mathbf{p})} = \frac{\text{Diag}(\mathbf{a}) - \mathbf{a} \mathbf{a}^\top}{\sqrt{d_q}}.$$

If \mathbf{a} approaches to a unit vector \mathbf{e} , then the Jacobian matrix $\frac{\partial \text{vec}(\mathbf{a})}{\partial \text{vec}(\mathbf{p})}$ will tend to $\mathbf{0}$.

Single-head self-attention

▷ Backward Process

$$\begin{aligned}\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{X})}, \\ &= (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\partial \text{vec}(\mathbf{A})}{\partial \text{vec}(\mathbf{P})} \frac{\partial \text{vec}(\mathbf{P})}{\partial \text{vec}(\mathbf{X})}.\end{aligned}$$

- the final Jacobian matrix

$$\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} \left((\mathbf{X}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{K} + (\mathbf{I}_n \otimes \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k) \right).$$

where

$$\mathbf{J} = \text{blockdiag}(\text{Diag}(\mathbf{A}_{:,1}) - \mathbf{A}_{:,1} \mathbf{A}_{:,1}^\top, \dots, \text{Diag}(\mathbf{A}_{:,n}) - \mathbf{A}_{:,n} \mathbf{A}_{:,n}^\top).$$

Single-head self-attention

Let us analysis each term

$$\frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})} = (\mathbf{A}^\top \otimes \mathbf{W}_v) + (\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} \left((\mathbf{X}^\top \mathbf{W}_k^\top \mathbf{W}_q \otimes \mathbf{I}_n) \mathbf{K} + (\mathbf{I}_n \otimes \mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k) \right).$$

What we see from the equation

- we have $\frac{\mathbf{J}}{\sqrt{d_q}}$
- we have $\mathbf{W}_q^\top \mathbf{W}_k$
- we even have $(\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n)$

Question:

what are the singular values of $\sigma \left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \right)$?

Single-head self-attention

Let us analysis each term in $\sigma \left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \right)$

- $\mathbf{W}_q^\top \mathbf{W}_k \in \mathcal{R}^{d \times d}$
- $(\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \in \mathcal{R}^{nd_v \times nn}$
- $\mathbf{J} \in \mathcal{R}^{nn \times nn}$
- $(\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \in \mathcal{R}^{nn \times nd}$
- $\left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \right) \in \mathcal{R}^{nd_v \times nd}$

Questions:

1. $\sigma(\mathbf{X} \mathbf{X}^\top)$
2. $\sigma(\mathbf{W}_q^\top \mathbf{W}_k)$
3. $\sigma(\mathbf{W}_v)$

Single-head self-attention

Our answers to some important questions in Transformer:

- 1. why Transformer is hard to train? NAN or INF
- 2. why normalization, \mathbf{X} , is important in Transformer?
- 3. why prenorm is more stable than postnorm? the norm of \mathbf{X} .
- 4. why deeper layer, L , is harder to train?
- 5. why larger hidden dimension, d , is harder to train?
- 6. why longer context, n , is harder to train?
- 7. why $\sigma(\mathbf{W}_q^\top \mathbf{W}_k)$ is important for model crash?
- 8. why learning rate warmup is important for classical Transformer?
- 9. why backward process always trigger problem in Transformer?
- 10. do we have solutions to so many why? Yes
- more...

1. why Transformer is hard to train?

Reason:

$$\sigma \left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \right)$$
 is prone to be very large.

2. why normalization is important in Transformer?

Reason:

because after normalization on \mathbf{X} , we can always promise $\|\mathbf{x}\|_2 = \sqrt{d}$ before γ, β , thus it will restrict $\sigma \left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n) \right)$ to be very large.

3. why prenorm is more stable than postnorm in Transformer?

Reason:

because in postnorm, $\mathbf{Y} = \mathbf{X} + \text{LN}(\text{sa}(\mathbf{X}))$, the norm of \mathbf{y} that is a column in \mathbf{Y} is always larger than \mathbf{x} (because in a high-dimensional space,

two vectors are prone to orthogonal to each other), thus, with along the

forward process, the norm of \mathbf{x} inputted into a self-attention module increases along with the increase of layer, in the backward process, it will

$$\text{trigger larger } \sigma \left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^T \mathbf{W}_q^T \mathbf{W}_k \otimes \mathbf{I}_n) \right).$$

- Derive the joint Jacobian matrix of LN + SA?
- Derive the joint Jacobian matrix of SA + LN?

4. why deeper layer is harder to train?

Reason:

because deeper layer means we will have more multiples of Jacobian matrix in chain rule, the backward process is more prone to have larger gradients.

5. why larger hidden dimension is harder to train?

Reason:

because larger d will be prone to lead to larger $\sigma(\mathbf{X}\mathbf{X}^\top)$ and larger $\mathbf{W}_q^\top \mathbf{W}_k$.

6. why longer context is harder to train?

Reason:

because larger n will be prone to lead to larger $\sigma(\mathbf{X}\mathbf{X}^\top)$.

7. why $\sigma(\mathbf{W}_q^\top \mathbf{W}_k)$ is important for model crash?

Reason:

because larger $\sigma(\mathbf{W}_q^\top \mathbf{W}_k)$ will be prone to lead to a larger
 $\sigma\left((\mathbf{I}_n \otimes \mathbf{W}_v \mathbf{X}) \frac{\mathbf{J}}{\sqrt{d_q}} (\mathbf{X}^\top \mathbf{W}_q^\top \mathbf{W}_k \otimes \mathbf{I}_n)\right).$

8. why learning rate warmup is important for classical Transformer?

Reason:

because learning rate warmup does not smooth $\mathbf{W}_q^\top \mathbf{W}_k$, but also smooth \mathbf{x} .

9. why backward process always triggers problem in classical Transformer

Reason:

because backward process is prone to lead to larger range of value, but in the forward process, we have normalization to constraint the values.

10. do we have solutions to so many why

Answer:

- constraint the singular values of \mathbf{X} ,
- constraint the singular values of $\mathbf{W}_q^\top \mathbf{W}_k$

Simulate gradient of self-attention

```
import torch
import torch.nn.functional as F
import numpy as np

def self_attention_with_grad(X, d_model=512):
    """
    Compute self-attention and its gradient with respect to input X
    Args:
        X: Input tensor of shape [batch_size, seq_len, d_model]
        d_model: Dimension of the model
    Returns:
        Y: Output tensor
        dy_dx: Gradient of Y with respect to X
    """
    # Ensure X requires gradient and create a fresh copy
    X = X.clone().detach().requires_grad_(True)

    batch_size, seq_len, _ = X.shape

    # Initialize learnable matrices
    W_q = torch.randn(d_model, d_model, device=X.device)*10
    W_k = torch.randn(d_model, d_model, device=X.device)*10
    W_v = torch.randn(d_model, d_model, device=X.device)

    # Compute Q, K, V
    Q = torch.matmul(X, W_q) # (batch_size, seq_len, d_model)
    K = torch.matmul(X, W_k) # (batch_size, seq_len, d_model)
    V = torch.matmul(X, W_v) # (batch_size, seq_len, d_model)

    # Scaled dot-product attention
    scale = torch.sqrt(torch.tensor(d_model, dtype=torch.float32, device=X.device))
    attention_scores = torch.matmul(Q, K.transpose(-2, -1)) / scale # (batch_size, seq_len, seq_len)
    attention_probs = F.softmax(attention_scores, dim=-1)

    # Compute output
    Y = torch.matmul(attention_probs, V) # (batch_size, seq_len, d_model)

    # Create a random gradient tensor to backpropagate
    grad_output = torch.randn_like(Y)

    # Compute gradient
    Y.backward(grad_output)

    # Get the gradient and make sure it exists
    assert X.grad is not None, "Gradient computation failed"
    grad = X.grad.clone()

    return Y.detach(), grad
```

Figure: Simulate gradient of self-attention

```
def main():
    # Set random seed for reproducibility
    torch.manual_seed(42)

    # Create sample input
    batch_size = 2
    seq_len = 1000
    d_model = 768
    X = torch.randn(batch_size, seq_len, d_model)

    # Compute self-attention and its gradient
    Y, X_gradient = self_attention_with_grad(X, d_model=d_model)

    # Print norms for verification
    print(f"Output Y norm: {Y.norm().item():.4f}")
    print(f"Gradient dY/dX norm: {X_gradient.norm().item():.4f}")

    # Print shapes for verification
    print(f"\nShapes:")
    print(f"Input X: {X.shape}")
    print(f"Output Y: {Y.shape}")
    print(f"Gradient dY/dX: {X_gradient.shape}")

    if __name__ == "__main__":
        main()

→ Output Y norm: 34595.2461
Gradient dY/dX norm: 42357236.0000

Shapes:
Input X: torch.Size([2, 1000, 768])
Output Y: torch.Size([2, 1000, 768])
Gradient dY/dX: torch.Size([2, 1000, 768])
```

Figure: Simulate gradient of self-attention.

Simulate gradient of FFN

```
import torch
import torch.nn.functional as F
import numpy as np

def ffn_with_grad(X, d_model=512, d_ff=2048):
    """
    Compute FFN (Feed-Forward Network) and its gradient with respect to input X
    Args:
        X: Input tensor of shape (batch_size, seq_len, d_model)
        d_model: Model dimension
        d_ff: Feed-forward dimension (usually 4*d_model)
    Returns:
        Y: Output tensor
        dL_dX: Gradient of L with respect to X
    """
    # Ensure X requires gradient and create a fresh copy
    X = X.clone().detach().requires_grad_(True)

    # Initialize weights for the two linear transformations
    W1 = torch.randn(d_model, d_ff, device=X.device) #/ np.sqrt(d_model)
    b1 = torch.zeros(d_ff, device=X.device)

    W2 = torch.randn(d_ff, d_model, device=X.device) #/ np.sqrt(d_ff)
    b2 = torch.zeros(d_model, device=X.device)

    # First linear transformation and ReLU
    hidden = torch.matmul(X, W1) + b1 # (batch_size, seq_len, d_ff)
    hidden = F.relu(hidden)

    # Second linear transformation
    Y = torch.matmul(hidden, W2) + b2 # (batch_size, seq_len, d_model)

    # Create a random gradient tensor to backpropagate
    grad_output = torch.randn_like(Y)
    print("Input grad_output norm: {:.4f}".format(grad_output.norm()))

    # Compute gradient
    Y.backward(grad_output)

    # Get the gradient and make sure it exists
    assert X.grad is not None, "Gradient computation failed"
    grad = X.grad.clone()

    return Y.detach(), grad
```

```
def main():
    # Set random seed for reproducibility
    torch.manual_seed(42)

    # Create sample input
    batch_size = 2
    seq_len = 1000
    d_model = 768
    d_ff = 768*4 # Usually 4*d_model

    X = torch.randn(batch_size, seq_len, d_model)

    # Compute FFN and its gradient
    Y, X_gradient = ffn_with_grad(X, d_model=d_model, d_ff=d_ff)

    # Print norms for verification
    print("Input X norm: {:.4f}".format(X.norm()))
    print("Output Y norm: {:.4f}".format(Y.norm()))
    print("Gradient dY/dX norm: {:.4f}".format(X_gradient.norm()))

    # Verify gradient shape matches input shape
    assert X.shape == X_gradient.shape, "Gradient shape mismatch!"

    if __name__ == "__main__":
        main()

    Input grad_output norm: 1239.5171
    Input X norm: 1239.5165
    Output Y norm: 1342592.1250
    Gradient dY/dX norm: 1346009.6250
```

Figure: Simulate gradient of FFN.

Figure: Simulate gradient of FFN.

Reproduce NaN problem

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class MultiHeadModel(nn.Module):
    def __init__(self, hidden_dim, num_heads):
        super(MultiHeadModel, self).__init__()
        self.hidden_dim = hidden_dim
        self.num_heads = num_heads
        self.head_dim = hidden_dim // num_heads
        scale = 1000

        rank = 2
        self.rank = rank

        # Low-rank decomposition for attention weights
        self.Aq = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(hidden_dim + rank, dtype=torch.float32)))
        self.Qb = nn.Parameter(scale * torch.randn(rank, hidden_dim) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qk = nn.Parameter(scale * torch.randn(rank, hidden_dim) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qv = nn.Parameter(scale * torch.randn(rank, hidden_dim) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qr = nn.Parameter(scale * torch.randn(rank, hidden_dim) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qz = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(hidden_dim + rank, dtype=torch.float32)))
        self.Bq = nn.Parameter(scale * torch.randn(rank, hidden_dim) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qa = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(hidden_dim + rank, dtype=torch.float32)))
        self.Qb = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qc = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qd = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qe = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qf = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qg = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))
        self.Qh = nn.Parameter(scale * torch.randn(hidden_dim, rank) / torch.sqrt(torch.tensor(rank + hidden_dim, dtype=torch.float32)))

        def scaled_dot_product_attention(self, Q, K, V):
            scores = Q.bmm(Q.transpose(-2, -1)) / (self.head_dim ** 0.5)
            attn_weights = F.softmax(scores, dim=-1)
            attn_output = torch.matmul(attn_weights, V)
            return attn_output

        def forward(self, x):
            # Compute low-rank Q, K, V
            Q = torch.matmul(torch.matmul(Q, self.Aq), self.Qb) # Low-rank Q
            K = torch.matmul(torch.matmul(Q, self.Qk), self.Qc) # Low-rank K
            V = torch.matmul(torch.matmul(Q, self.Qv), self.Qd) # Low-rank V

            # Reshape for multi-head attention
            Q = Q.view(Q.shape[0], Q.shape[1], self.num_heads, self.head_dim).permute(0, 1, 3)
            K = K.view(K.shape[0], K.shape[1], self.num_heads, self.head_dim).permute(0, 1, 3)
            V = V.view(V.shape[0], V.shape[1], self.num_heads, self.head_dim).permute(0, 1, 3)

            # Scaled dot-product attention
            attn_output = self.scaled_dot_product_attention(Q, K, V)
            attn_output = attn_output.permute(0, 1, 3).contiguous()
            attn_output = attn_output.view(attn_output.shape[0], attn_output.shape[1], -1)

            # Output projection
            attn_output = torch.matmul(attn_output, self.Qf)
            return attn_output
```

```
+ Code + Text
✓  class FeedForwardNetwork(nn.Module):
    def __init__(self, hidden_dim, ff_dim):
        super(FeedForwardNetwork, self).__init__()
        scale = 1.0
        self.W1 = nn.Parameter(scale * torch.randn(hidden_dim, ff_dim) / torch.sqrt(torch.tensor(hidden_dim + ff_dim, dtype=torch.float32)))
        self.W2 = nn.Parameter(scale * torch.randn(ff_dim, hidden_dim) / torch.sqrt(torch.tensor(hidden_dim + ff_dim, dtype=torch.float32)))
        self.B1 = nn.Parameter(torch.zeros(ff_dim))

    def forward(self, x):
        ffn_output = torch.matmul(x, self.W1)
        ffn_output = F.relu(ffn_output)
        ffn_output = torch.matmul(ffn_output, self.W2)
        return ffn_output

class PreNormTransformerBlock(nn.Module):
    def __init__(self, hidden_dim, num_heads, ff_dim):
        super(PreNormTransformerBlock, self).__init__()
        self.attention = MultiHeadModel(hidden_dim, num_heads)
        self.ffn = FeedForwardNetwork(hidden_dim, ff_dim)
        self.ln1 = nn.LayerNorm(hidden_dim)
        self.ln2 = nn.LayerNorm(hidden_dim)

    def forward(self, x):
        # Pre-norm and attention
        x1 = x
        x1 = self.ln1(x1)
        x1 = self.attention(x1)
        x1_norm = x1.norm(2) # Retain gradients for input to attention
        attn_output = self.attention(x1_norm)
        x2 = x1 + attn_output

        # Pre-norm and feed-forward network
        x2_norm = x2.norm(2)
        x2_norm.retain_grad() # Retain gradients for input to FFN
        ffn_output = self.ffn(x2_norm)
        x = x2 + ffn_output
        return x, x1, x2, x1_norm, x2_norm

class PreNormTransformerModel(nn.Module):
    def __init__(self, num_layers, hidden_dim, num_heads, ff_dim):
        super(PreNormTransformerModel, self).__init__()
        self.layers = nn.ModuleList([PreNormTransformerBlock(hidden_dim, num_heads, ff_dim) for _ in range(num_layers)])
        self.attention_inputs = []
        self.ffn_inputs = []
        self.ln1_norms = []

    def forward(self, x):
        self.attention_inputs = []
        self.ffn_inputs = []
        for i, layer in enumerate(self.layers):
            x1, x2, x1_norm, x2_norm = layer(x)

            norm1 = torch.norm(x1, dim=1).mean() # Calculate the mean norm of activations
            print(f"Norm input into self-attention block {i+1}: {norm1.item():.4f}")

            norm2 = torch.norm(x2, dim=1).mean() # Calculate the mean norm of activations
            print(f"Norm input into FFN block {i+1}: {norm2.item():.4f}")

            self.attention_inputs.append(x1_norm)
            self.ffn_inputs.append(x2_norm)
            self.ln1_norms.append(x1_norm)

        return x
```

Figure: Reproduce NaN problem (part 1).

Figure: Reproduce NaN problem (part 2).

Reproduce NaN problem

+ Code + Text

```
✓ 0 # Define input tensor and model parameters
batch_size = 1
seq_len = 1000
hidden_dim = 1024
num_layers = 12
num_heads = 8
ff_dim = 4096

# Instantiate model
model = PreformerTransformerModel(num_layers, hidden_dim, num_heads, ff_dim)

# Random input tensor
X = torch.rand(batch_size, seq_len, hidden_dim, requires_grad=True)

device = torch.device("cuda") if torch.cuda.is_available() else "cpu"
model = model.to(device)
X = X.to(device)

# Define a dummy loss function
dummy_target = torch.ones(batch_size, seq_len, hidden_dim).to(device)
criterion = nn.MSELoss()

# Forward pass
output = model(X)

# Compute loss
loss = criterion(output, dummy_target)

# Backward pass
loss.backward()

# Output gradients for inputs to attention and FFN blocks
for i, (atten_input, ffn_input) in enumerate(zip(model.attention_inputs, model.ffn_inputs)):
    print(f"Gradient of L w.r.t. input to attention block [{i+1}]: {torch.mean(atten_input.grad.norm(dim=2)).item():.4f}")
    print(f"Gradient of L w.r.t. input to FFN block [{i+1}]: {torch.mean(ffn_input.grad.norm(dim=2)).item():.4f}")
```

Figure: Reproduce NaN problem (part 3).

```
✓ Norm input into self-attention block 1: 31.9515
Norm input into FFN block 1: 1504378304.0000
Norm input into self-attention block 2: 1504378304.0000
Norm input into FFN block 2: 1796688128.0000
Norm input into self-attention block 3: 1796688128.0000
Norm input into FFN block 3: 2342827264.0000
Norm input into self-attention block 4: 2342827264.0000
Norm input into FFN block 4: 2449437440.0000
Norm input into self-attention block 5: 2449437440.0000
Norm input into FFN block 5: 2779644160.0000
Norm input into self-attention block 6: 2779644160.0000
Norm input into FFN block 6: 3114456576.0000
Norm input into self-attention block 7: 3114456576.0000
Norm input into FFN block 7: 3297453312.0000
Norm input into self-attention block 8: 3297453312.0000
Norm input into FFN block 8: 3485641472.0000
Norm input into self-attention block 9: 3485641472.0000
Norm input into FFN block 9: 3745204224.0000
Norm input into self-attention block 10: 3745204224.0000
Norm input into FFN block 10: 3945027072.0000
Norm input into self-attention block 11: 3945027072.0000
Norm input into FFN block 11: 4139921216.0000
Norm input into self-attention block 12: 4139921216.0000
Norm input into FFN block 12: 4241150976.0000
Gradient of L w.r.t. input to attention block 1: nan
Gradient of L w.r.t. input to FFN block 1: nan
Gradient of L w.r.t. input to attention block 2: inf
Gradient of L w.r.t. input to FFN block 2: inf
Gradient of L w.r.t. input to attention block 3: inf
Gradient of L w.r.t. input to FFN block 3: 7907124449378384.0000
Gradient of L w.r.t. input to attention block 4: inf
Gradient of L w.r.t. input to FFN block 4: 15931159411712.0000
Gradient of L w.r.t. input to attention block 5: inf
Gradient of L w.r.t. input to FFN block 5: 438227604608.0000
Gradient of L w.r.t. input to attention block 6: inf
Gradient of L w.r.t. input to FFN block 6: 8214941824.0000
Gradient of L w.r.t. input to attention block 7: inf
Gradient of L w.r.t. input to FFN block 7: 1358573856.0000
Gradient of L w.r.t. input to attention block 8: inf
Gradient of L w.r.t. input to FFN block 8: 47305944.0000
Gradient of L w.r.t. input to attention block 9: inf
Gradient of L w.r.t. input to FFN block 9: 2325342.2500
Gradient of L w.r.t. input to attention block 10: 983398400888031232.0000
Gradient of L w.r.t. input to FFN block 10: 663090.0075
Gradient of L w.r.t. input to attention block 11: 2073521025843200.0000
Gradient of L w.r.t. input to FFN block 11: 1942.0000
Gradient of L w.r.t. input to attention block 12: 9043992287744.0000
Gradient of L w.r.t. input to FFN block 12: 23.7339
```

Figure: Reproduce NaN problem (part 4).

RMSNorm

- Forward Process

$$\text{RMSN}(\mathbf{x}) = \gamma \odot \frac{\sqrt{d}\mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}$$

where $\mathbf{x} \in \mathcal{R}^d$

- Backward Process

$$\frac{\partial \text{RMSN}(\mathbf{x})}{\partial \mathbf{x}} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|_2^2 + \epsilon} \right) \text{Diag}(\gamma).$$

RMSNorm gradient

Matrix Calculus

MatrixCalculus provides matrix calculus for everyone. It is an online tool that computes vector and matrix derivatives (matrix calculus).

The screenshot shows the MatrixCalculus interface. In the input field, the expression $r \cdot \mathbf{x} / (\mathbf{x}' \cdot \mathbf{x} + e)^{0.5}$ is entered, with "w.r.t." set to \mathbf{x} . The resulting derivative is displayed as:

$$\frac{\partial}{\partial \mathbf{x}} \left((r \odot \mathbf{x}) / (\mathbf{x}^\top \cdot \mathbf{x} + e)^{0.5} \right) = \\ \frac{1}{t_0^{0.5}} \cdot \text{diag}(r) - \frac{1}{t_0^{1.5}} \cdot r \odot \mathbf{x} \cdot \mathbf{x}^\top$$

where

$$t_0 = e + \mathbf{x}^\top \cdot \mathbf{x}$$

Below the input field, there are dropdown menus for variable types: "e is a scalar", "r is a vector", and "x is a vector". To the right, there are buttons for "Export functions as Python" (which is selected), "Latex", and "Common subexpressions ON".

Figure: RMSNorm gradient.

StableNorm

- Forward Process

$$\text{StableNorm}(\mathbf{x}) = \gamma \odot \frac{d^\alpha \mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}$$

where $\mathbf{x} \in \mathcal{R}^d$

- Backward Process

$$\frac{\partial \text{StableNorm}(\mathbf{x})}{\partial \mathbf{x}} = \frac{d^\alpha}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{x} \mathbf{x}^\top}{\|\mathbf{x}\|_2^2 + \epsilon} \right) \text{Diag}(\gamma).$$

L_2 Norm

- Forward Process

$$\text{L2Norm}(\mathbf{x}) = \frac{\mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}}$$

where $\mathbf{x} \in \mathcal{R}^d$

- Backward Process

$$\frac{\partial \text{L2Norm}(\mathbf{x})}{\partial \mathbf{x}} = \frac{1}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|_2^2 + \epsilon} \right).$$

Self-attention with QKNorm

- Forward process of self-attention with QKNorm

$$\mathbf{A} = \text{softmax}(\mathbf{P}^{(2)}), \quad \text{where } P_{ij}^{(2)} = \left(\frac{\text{RMSN}(\mathbf{W}_q \mathbf{x}_i)^\top \text{RMSN}(\mathbf{W}_k \mathbf{x}_j)}{\sqrt{d_1}} \right).$$

- Backward Process

$$\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_i} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \mathbf{W}_q^\top \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{Diag}(\boldsymbol{\gamma}_q) \text{Diag}(\boldsymbol{\gamma}_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}},$$

$$\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{x}_j} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \mathbf{W}_k^\top \left(\mathbf{I} - \frac{\mathbf{W}_k \mathbf{x}_j (\mathbf{W}_k \mathbf{x}_j)^\top}{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon} \right) \text{Diag}(\boldsymbol{\gamma}_k) \text{Diag}(\boldsymbol{\gamma}_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}},$$

$$\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_q} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{W}_q \mathbf{x}_i (\mathbf{W}_q \mathbf{x}_i)^\top}{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon} \right) \text{Diag}(\boldsymbol{\gamma}_q) \text{Diag}(\boldsymbol{\gamma}_k) \frac{\mathbf{W}_k \mathbf{x}_j}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \mathbf{x}_i^\top,$$

$$\frac{\partial P_{ij}^{(2)}}{\partial \mathbf{W}_k} = \frac{\sqrt{d}}{\sqrt{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon}} \left(\mathbf{I} - \frac{\mathbf{W}_k \mathbf{x}_j (\mathbf{W}_k \mathbf{x}_j)^\top}{\|\mathbf{W}_k \mathbf{x}_j\|_2^2 + \epsilon} \right) \text{Diag}(\boldsymbol{\gamma}_k) \text{Diag}(\boldsymbol{\gamma}_q) \frac{\mathbf{W}_q \mathbf{x}_i}{\sqrt{\|\mathbf{W}_q \mathbf{x}_i\|_2^2 + \epsilon}} \mathbf{x}_j^\top.$$

Feed-Forward Network with SwishGLU

▷ Forward Process

$$\mathbf{u} = \mathbf{W}_1 \mathbf{x}, \mathbf{v} = \mathbf{W}_2 \mathbf{x}, \mathbf{z} = \mathbf{u} \odot \text{SiLU}(\mathbf{v}),$$

$$\mathbf{y} = \mathbf{W}_3 \mathbf{z}.$$

where $\text{SiLU}(\mathbf{v}) = \mathbf{v} \odot \sigma(\mathbf{v})$

▷ Backward Process

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{W}_3 [\text{Diag}(\text{SiLU}(\mathbf{W}_2 \mathbf{x})) \mathbf{W}_1 + \text{Diag}(\mathbf{W}_1 \mathbf{x}) \text{Diag}(\text{SiLU}'(\mathbf{W}_2 \mathbf{x})) \mathbf{W}_2]$$

where $\text{Diag}(\text{SiLU}'(\mathbf{v})) = \text{Diag}(\sigma(\mathbf{v}) + \mathbf{v} \sigma'(\mathbf{v})(1 - \sigma(\mathbf{v}))).$

Five Types of Normalizations

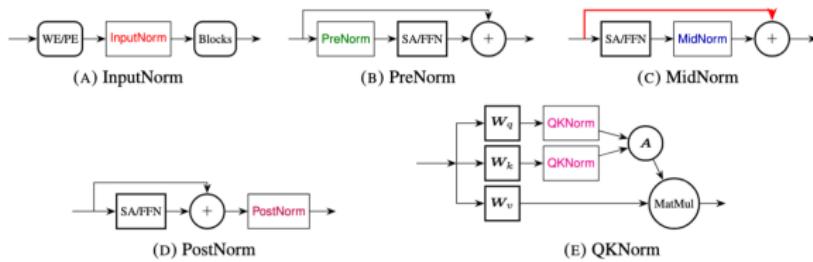


FIGURE 2: Five different normalization methods. The only difference between them is the position of normalization. In (A), WE/PE indicates word embedding and patch embedding.

Figure: Five types of normalizations according to their position in the network.

Deeply Normalized Transformer (DNT)

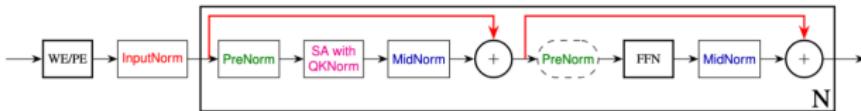


FIGURE 3: DNT architecture. The second PreNorm marked with dashed and rounded corners is optional. By default, we do not use the second PreNorm.

Figure: DNT network.

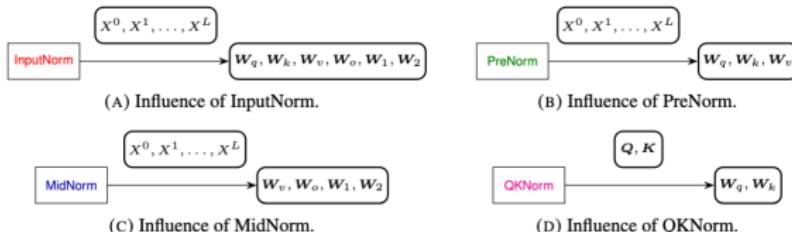


FIGURE 4: Influence of different normalizations. For instance, InputNorm stabilizes $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o, \mathbf{W}_1, \mathbf{W}_2$ by constraining X^0, X^1, \dots, X^L .

Figure: Value analysis of each normalization in DNT.

Unknown problems of mine

We still have many unknown problems

- how to characterize the representation capabilities of transformers?
Lipschitz constant? smoothness?
- how to store memory in Transformer?
- how to assess one model without test set?
- what are the characters of a stronger model?
- how can we invent an optimizer that largely improves Adam?
- more and more...

References I

- [1] Sheldon Axler. *Linear algebra done right*. Springer, 2015.
- [2] Christopher M Bishop and Hugh Bishop. *Deep learning: Foundations and concepts*. Springer Nature, 2023.
- [3] Stanley H Chan. *Introduction to probability for data science*. Michigan Publishing, 2021.
- [4] Jeffrey R Chasnov. *Differential Equations for Engineers*. 2019.
- [5] Jeffrey R Chasnov. *Matrix Algebra for Engineers*. 2018.
- [6] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge University Press, 2020.
- [7] Alan Edelman and Steven G Johnson. “18. S096 Matrix Calculus for Machine Learning and Beyond, January IAP 2022”. In: (2022).

References II

- [8] Lawrence C Evans. *An introduction to stochastic differential equations*. Vol. 82. American Mathematical Soc., 2012.
- [9] Alexander Graham. *Kronecker products and matrix calculus with applications*. Courier Dover Publications, 2018.
- [10] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- [11] Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge university press, 1994.
- [12] Hyunjik Kim, George Papamakarios, and Andriy Mnih. “The lipschitz constant of self-attention”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 5562–5571.
- [13] Kevin P Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.

References III

- [14] Yurii Nesterov. "Introductory lectures on convex programming volume i: Basic course". In: *Lecture notes 3.4* (1998), p. 5.
- [15] Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999.
- [16] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. "The matrix cookbook". In: *Technical University of Denmark 7.15* (2008), p. 510.
- [17] Xianbiao Qi, Jianan Wang, and Lei Zhang. "Understanding optimization of deep learning via jacobian matrix and lipschitz constant". In: *arXiv preprint arXiv:2306.09338* (2023).
- [18] Xianbiao Qi et al. "Lipsformer: Introducing lipschitz continuity to vision transformers". In: *arXiv preprint arXiv:2304.09856* (2023).
- [19] Xianbiao Qi et al. "Stable-Transformer: towards a stable transformer training". In: *arXiv* (2025).

References IV

- [20] Xianbiao Qi et al. "Taming transformer without using learning rate warmup". In: *ICLR* (2025).
- [21] Roland Speicher. *High-Dimensional Analysis: Random Matrices and Machine Learning*. Lecture Notes. Department of Mathematics, Saarland University, 2023.
- [22] Gilbert Strang. *Linear algebra and learning from data*. SIAM, 2019.
- [23] Gilbert Strang, Kenji Hiranabe, and Ashley Fernandes. "The Art of Linear Algebra". In: *PRIMUS* (2024), pp. 1–14.
- [24] Lloyd N Trefethen and David Bau. *Numerical linear algebra*. SIAM, 2022.
- [25] A Vaswani. "Attention is all you need". In: *Advances in Neural Information Processing Systems* (2017).

References V

- [26] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*. Vol. 47. Cambridge university press, 2018.
- [27] Ruibin Xiong et al. “On layer normalization in the transformer architecture”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 10524–10533.
- [28] Raymond W Yeung. *A first course in information theory*. Springer Science & Business Media, 2012.
- [29] Biao Zhang and Rico Sennrich. “Root mean square layer normalization”. In: *Advances in Neural Information Processing Systems 32* (2019).
- [30] Xianda Zhang. *Matrix analysis and applications*. Cambridge University Press, 2017.