

Chapter7 Presentation

ZhouXuRui

HITSZ

December 9, 2022

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Vector

A vector $\mathbf{x} \in \mathbb{R}^n$ is a list of n numbers, usually written as a column vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

The unit vector \mathbf{e}_i is a vector of all 0's, except entry i , which has value 1:

$$\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)$$

This is also called a one-hot vector.

Matrix

A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with m rows and n columns is a 2 d array of numbers, arranged as follows:

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

If $m = n$, the matrix is said to be square.

We use the notation A_{ij} or $A_{i,j}$ to denote the entry of \mathbf{A} in the i th row and j th column. We use the notation $\mathbf{A}_{i,:}$ to denote the i 'th row and $\mathbf{A}_{:,j}$ to denote the j 'th column. We treat all vectors as column vectors by default (so $\mathbf{A}_{i,:}$ is viewed as a column vector with n entries). We use bold upper case letters to denote matrices, bold lower case letters to denote vectors, and non-bold letters to denote scalars.

Notation

We can view a matrix as a set of columns stacked along the horizontal axis:

$$\mathbf{A} = \left[\begin{array}{c|c|c|c} \mathbf{A}_{:,1} & \mathbf{A}_{:,2} & \cdots & \mathbf{A}_{:,n} \end{array} \right]$$

For brevity, we will denote this by

$$\mathbf{A} = [\mathbf{A}_{:,1}, \mathbf{A}_{:,2}, \dots, \mathbf{A}_{:,n}]$$

We can also view a matrix as a set of rows stacked along the vertical axis:

$$\mathbf{A} = \left[\begin{array}{c} - \\ - \\ \vdots \\ - \end{array} \begin{array}{c} \mathbf{A}_{1,:}^\top \\ \mathbf{A}_{2,:}^\top \\ \vdots \\ \mathbf{A}_{m,:}^\top \end{array} \begin{array}{c} - \\ - \\ \vdots \\ - \end{array} \right]$$

For brevity, we will denote this by

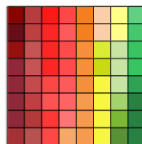
$$\mathbf{A} = [\mathbf{A}_{1,:}; \mathbf{A}_{2,:}; \dots; \mathbf{A}_{m,:}]$$

Tensor

A tensor (in machine learning terminology) is just a generalization of a 2d array to more than 2 dimensions. The number of dimensions is known as the order or rank of the tensor.

In mathematics, tensors can be viewed as a way to define multilinear maps, just as matrices can be used to define linear functions, although we will not need to use this interpretation.

Vector



$\mathbb{R}^{8 \times 8}$



$\mathbb{R}^{4 \times 4 \times 4}$

Notation

More introduction about Tensor:

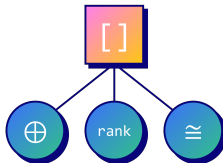
Tensors in machine learning and Math/Physics:



Data structure

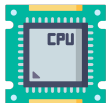


Object



Tensors in deep learning:

Arrays



`np.array`

vs.

Tensors



`torch.tensor`

To put it more clearly:

Machine learning Tensor

In machine learning, the term “rank r tensor” is just another way to say “an r dimensional array of numbers.” It just a way to refer to a data structure. There is no mathematical content to the usage of the term.

Mathematics Tensor

The tensors of mathematics are the linear representations of multilinear maps, and vector bundles constructed from them. The ideas make sense for arbitrary modules over arbitrary rings.

Pyhsics Tensor

The tensors of physics are sections of vector bundles on manifolds modeling configuration spaces/spacetime which are themselves tensor products of tangent and cotangent spaces. This means that relative to a choice of coordinate systems on the manifold, they are functions taking values in multilinear maps whose arguments are either vectors like the partial derivative operators in each of the coordinate directions, or forms like the differentials of the coordinate functions. Evaluating them on these basis elements yields a multidimensional array of functions.

1 Section 1: Introduction

- Notation
- **Vector Space**
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Vector Space

Space in Mathematics

A space consists of selected mathematical objects that are treated as points, and selected relationships(structure) between these points.

Vector Space

Definition: A vector space consists of a set V (elements of V are called vectors), a field \mathbb{F} (elements of \mathbb{F} are called scalars), and two operations

- An operation called vector addition that takes two vectors $v, w \in V$, and produces a third vector, written $v + w \in V$.
- An operation called scalar multiplication that takes a scalar $c \in \mathbb{F}$ and a vector $v \in V$, and produces a new vector, written $cv \in V$ which satisfy the following conditions (called axioms).

Axioms

1. Associativity of vector addition: $(u + v) + w = u + (v + w)$ for all $u, v, w \in V$.
2. Existence of a zero vector: There is a vector in V , written 0 and called the zero vector, which has the property that $u + 0 = u$ for all $u \in V$.
3. Existence of negatives: For every $u \in V$, there is a vector in V , written $-u$ and called the negative of u , which has the property that $u + (-u) = 0$.
4. Associativity of multiplication: $(ab)u = a(bu)$ for any $a, b \in \mathbb{F}$ and $u \in V$.
5. Distributivity: $(a + b)u = au + bu$ and $a(u + v) = au + av$ for all $a, b \in \mathbb{F}$ and $u, v \in V$.
6. Unitality: $1u = u$ for all $u \in V$.

Vector Space

Linear Independence

A set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is said to be (linearly) independent if no vector can be represented as a linear combination of the remaining vectors. Conversely, a vector which can be represented as a linear combination of the remaining vectors is said to be (linearly) dependent.

Span

The span of a set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is the set of all vectors that can be expressed as a linear combination of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. That is,

$$\text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\}) \triangleq \left\{ \mathbf{v} : \mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{x}_i, \quad \alpha_i \in \mathbb{R} \right\}$$

Basis

A basis \mathcal{B} is a set of linearly independent vectors that spans the whole space, meaning that $\text{span}(\mathcal{B}) = \mathbb{R}^n$.

There are often multiple bases to choose from, the standard basis uses the coordinate vectors $\mathbf{e}_1 = (1, 0, \dots, 0)$, up to $\mathbf{e}_n = (0, 0, \dots, 0, 1)$. This lets us translate back and forth between viewing a vector in \mathbb{R}^2 as either an "arrow in the plane", rooted at the origin, or as an ordered list of numbers (corresponding to the coefficients for each basis vector).

Linear Transformation

A linear map or linear transformation is any function $f : \mathcal{V} \rightarrow \mathcal{W}$ such that $f(\mathbf{v} + \mathbf{w}) = f(\mathbf{v}) + f(\mathbf{w})$ and $f(a\mathbf{v}) = af(\mathbf{v})$ for all $\mathbf{v}, \mathbf{w} \in \mathcal{V}$.

Once the basis of \mathcal{V} is chosen, a linear map $f : \mathcal{V} \rightarrow \mathcal{W}$ is completely determined by specifying the images of the basis vectors, because any element of \mathcal{V} can be expressed uniquely as a linear combination of them.

Suppose $\mathcal{V} = \mathbb{R}^n$ and $\mathcal{W} = \mathbb{R}^m$. We can compute $f(\mathbf{v}_i) \in \mathbb{R}^m$ for each basis vector in \mathcal{V} , and store these along the columns of an $m \times n$ matrix **A**. We can then compute $\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^m$ for any $\mathbf{x} \in \mathbb{R}^n$ as follows:

$$\mathbf{y} = \left(\sum_{j=1}^n a_{1j}x_j, \dots, \sum_{j=1}^n a_{mj}x_j \right)$$

This corresponds to multiplying the vector \mathbf{x} by the matrix \mathbf{A} :

$$\mathbf{y} = \mathbf{Ax}$$

That is, Matrix can be viewed as Linear Transformation.

Span

Suppose we view a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ as a set of n vectors in \mathbb{R}^m . The range (sometimes also called the column space) of this matrix is the span of the columns of \mathbf{A} . In other words,

$$\text{range}(\mathbf{A}) \triangleq \{\mathbf{v} \in \mathbb{R}^m : \mathbf{v} = \mathbf{A}\mathbf{x}, \mathbf{x} \in \mathbb{R}^n\}.$$

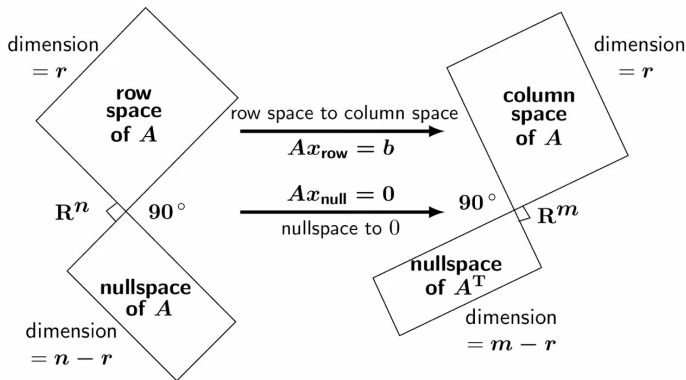
This can be thought of as the set of vectors that can be "reached" or "generated" by \mathbf{A} ; it is a subspace of \mathbb{R}^m whose dimensionality is given by the rank of \mathbf{A} . The nullspace of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the set of all vectors that get mapped to the null vector when multiplied by \mathbf{A} , i.e.,

$$\text{nullspace}(\mathbf{A}) \triangleq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

Vector Space

Here is the big picture of linear algebra in MIT 18.06:

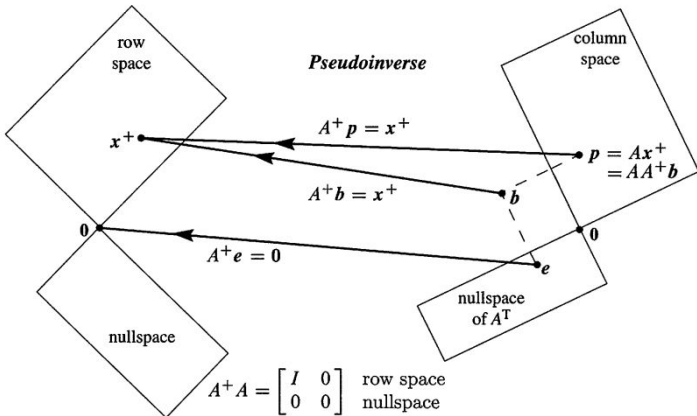
Big Picture of Linear Algebra



This is the Big Picture—two subspaces in \mathbb{R}^n and two subspaces in \mathbb{R}^m .

From row space to column space, A is invertible.

Vector Space



Projection

The projection of a vector $\mathbf{y} \in \mathbb{R}^m$ onto the span of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ (here we assume $\mathbf{x}_i \in \mathbb{R}^m$) is the vector $\mathbf{v} \in \text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\})$, such that \mathbf{v} is as close as possible to \mathbf{y} , as measured by the Euclidean norm $\|\mathbf{v} - \mathbf{y}\|_2$. We denote the projection as $\text{Proj}(\mathbf{y}; \{\mathbf{x}_1, \dots, \mathbf{x}_n\})$ and can define it formally as

$$\text{Proj}(\mathbf{y}; \{\mathbf{x}_1, \dots, \mathbf{x}_n\}) = \underset{\mathbf{v} \in \text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\})}{\text{argmin}} \|\mathbf{y} - \mathbf{v}\|_2$$

Given a (full rank) matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$, we can define the projection of a vector $\mathbf{y} \in \mathbb{R}^m$ onto the range of \mathbf{A} as follows:

$$\text{Proj}(\mathbf{y}; \mathbf{A}) = \underset{\mathbf{v} \in \mathcal{R}(\mathbf{A})}{\text{argmin}} \|\mathbf{v} - \mathbf{y}\|_2 = \mathbf{A} \left(\mathbf{A}^\top \mathbf{A} \right)^{-1} \mathbf{A}^\top \mathbf{y}.$$

Proof

When you want to find the projection p of b onto $C(A)$, then you can start with a "parameterization" of p that guarantees that your result is in $C(A)$. So you set $p = A\bar{x}$ and $\bar{x} \in V$. Now you want $p - b$ to be orthogonal to $C(A)$, which means it's in the null space of A^T , so we can get $A^T(p - b) = 0$ or $A^T(A\bar{x} - b) = 0$ or

$$\bar{x} = (A^T A)^{-1} A^T b$$

Now you have the particular $\bar{x} \in V$ that provides the correct parameters for your projection p , and you just have to apply your initial choice $p = A\bar{x}$ for the parametrization to obtain the projection p .

Vector Space

When $Ax = b$ is inconsistent, its least-squares solution minimizes $\|Ax - b\|^2$:

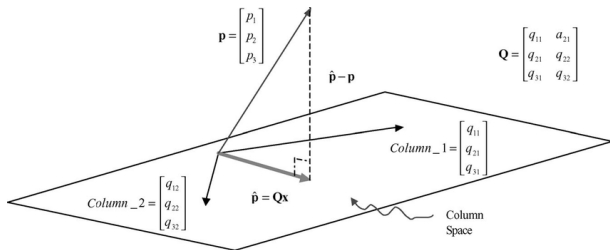
- Normal equations $A^T A \hat{x} = A^T b$.

$A^T A$ is invertible exactly when the columns of A are linearly independent! Then,

- Best estimate \hat{x} $\hat{x} = (A^T A)^{-1} A^T b$.

The projection of b onto the column space is the nearest point $A\hat{x}$:

- Projection $p = A\hat{x} = A(A^T A)^{-1} A^T b$.



1 Section 1: Introduction

- Notation
- Vector Space
- **Norms**
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Vector Norms

A norm of a vector $\|\mathbf{x}\|$ is, informally, a measure of the "length" of the vector. More formally, a norm is any function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ that satisfies 4 properties:

- For all $\mathbf{x} \in \mathbb{R}^n$, $f(\mathbf{x}) \geq 0$ (non-negativity).
- For all $\mathbf{x} \in \mathbb{R}^n$, $t \in \mathbb{R}$, $f(t\mathbf{x}) = |t|f(\mathbf{x})$ (absolute value homogeneity).
- For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (triangle inequality).

Consider the following common examples:

p-norm $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$, for $p \geq 1$

2-norm $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$, also called Euclidean norm. Note that

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x}.$$

1-norm $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$.

Max-norm $\|\mathbf{x}\|_\infty = \max_i |x_i|$.

Matrix Norms

Suppose we think of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ as defining a linear function $f(\mathbf{x}) = \mathbf{Ax}$. We define the induced norm of \mathbf{A} as the maximum amount by which f can lengthen any unit-norm input:

$$\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p} = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|_p$$

Typically $p = 2$, in which case

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} = \max_i \sigma_i$$

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- **Properties**
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Trace

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

$$\text{tr}(\mathbf{A}) \triangleq \sum_{i=1}^n A_{ii}.$$

The trace has the following properties, where $c \in \mathbb{R}$ is a scalar, and $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ are square matrices:

$$\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A}^\top)$$

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

$$\text{tr}(c\mathbf{A}) = c \text{tr}(\mathbf{A})$$

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

Properties

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i \text{ where } \lambda_i \text{ are the eigenvalues of } \mathbf{A}$$

We also have the following important cyclic permutation property: For $\mathbf{A}, \mathbf{B}, \mathbf{C}$ such that \mathbf{ABC} is square,

$$\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB})$$

From this, we can derive the trace trick, which rewrites the scalar inner product $\mathbf{x}^\top \mathbf{Ax}$ as follows

$$\mathbf{x}^\top \mathbf{Ax} = \text{tr}(\mathbf{x}^\top \mathbf{Ax}) = \text{tr}(\mathbf{xx}^\top \mathbf{A})$$

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- **Special Types**

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Positive Definite Matrix

Given a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and a vector $\mathbf{x} \in \mathbb{R}^n$, the scalar value $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is called a quadratic form. Written explicitly, we see that

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j.$$

Note that,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \left(\mathbf{x}^\top \mathbf{A} \mathbf{x} \right)^\top = \mathbf{x}^\top \mathbf{A}^\top \mathbf{x} = \mathbf{x}^\top \left(\frac{1}{2} \mathbf{A} + \frac{1}{2} \mathbf{A}^\top \right) \mathbf{x}$$

For this reason, we often implicitly assume that the matrices appearing in a quadratic form are symmetric.

Special Types

Positive Definite Matrix

The following conditions are equivalent:

- 1) $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \neq 0$.
- 2) All eigenvalues of A satisfy $\lambda_j > 0$.
- 3) All NW (upper left) minors of A are positive.
- 4) In the reduction of A to row echelon form, no exchanges are required, and all pivots are positive.
- 5) $A = R^T R$ for some non-singular R .

Positive Definite Matrix

For a symmetric matrix A , the following conditions are equivalent.

- (1) $A \succeq 0$.
- (2) $A = U^T U$ for some matrix U .
- (3) $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for every $\mathbf{x} \in \mathbb{R}^n$.
- (4) All principal minors of A are nonnegative.

AA^T is positive semidefinite.

Proof: Let $y = A^T x$.

$$x^T AA^T x = y^T y = \sum_{k=1}^N y_k^2 \geq 0$$

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- **Vector–vector products**
- Matrix–vector products
- Matrix–matrix products
- Manipulating data matrices

Vector–vector products

Given two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, the quantity $\mathbf{x}^\top \mathbf{y}$, called the inner product, dot product or scalar product of the vectors, is a real number given by

$$\langle \mathbf{x}, \mathbf{y} \rangle \triangleq \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i$$

Given vectors $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ (they no longer have to be the same size), $\mathbf{x}\mathbf{y}^\top$ is called the outer product of the vectors. It is a matrix whose entries are given by $(\mathbf{x}\mathbf{y}^\top)_{ij} = x_i y_j$, i.e.,

$$\mathbf{x}\mathbf{y}^\top \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}.$$

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- **Matrix–vector products**
- Matrix–matrix products
- Manipulating data matrices

Matrix–vector products

Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a vector $\mathbf{x} \in \mathbb{R}^n$, their product is a vector $\mathbf{y} = \mathbf{Ax} \in \mathbb{R}^m$. There are a couple ways of looking at matrix-vector multiplication, and we will look at them both.

If we write \mathbf{A} by rows, then we can express $\mathbf{y} = \mathbf{Ax}$ as follows:

$$\mathbf{y} = \mathbf{Ax} = \begin{bmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ & \vdots & \\ - & \mathbf{a}_m^\top & - \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_m^\top \mathbf{x} \end{bmatrix}$$

Matrix–vector products

In other words, the i th entry of \mathbf{y} is equal to the inner product of the i th row of \mathbf{A} and \mathbf{x} , $y_i = \mathbf{a}_i^\top \mathbf{x}$. Alternatively, let's write \mathbf{A} in column form. In this case we see that

$$\mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{a}_1 \\ | \end{bmatrix} x_1 + \cdots + \begin{bmatrix} | \\ \mathbf{a}_n \\ | \end{bmatrix} x_n.$$

In other words, \mathbf{y} is a linear combination of the columns of \mathbf{A} , where the coefficients of the linear combination are given by the entries of \mathbf{x} . We can view the columns of \mathbf{A} as a set of basis vectors defining a linear subspace. We can construct vectors in this subspace by taking linear combinations of the basis vectors.

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- **Matrix–matrix products**
- Manipulating data matrices

Matrix–matrix products

Below we look at four different (but, of course, equivalent) ways of viewing the matrix-matrix multiplication $\mathbf{C} = \mathbf{AB}$.

First we can view matrix-matrix multiplication as a set of vector-vector products. The most obvious viewpoint, which follows immediately from the definition, is that the i, j entry of \mathbf{C} is equal to the inner product of the i th row of \mathbf{A} and the j th column of \mathbf{B} . Symbolically, this looks like the following:

$$\begin{aligned}\mathbf{C} = \mathbf{AB} &= \begin{bmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ & \vdots & \\ - & \mathbf{a}_m^\top & - \end{bmatrix} \begin{bmatrix} | & | & \cdots & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & | & & | \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}_1^\top \mathbf{b}_1 & \mathbf{a}_1^\top \mathbf{b}_2 & \cdots & \mathbf{a}_1^\top \mathbf{b}_p \\ \mathbf{a}_2^\top \mathbf{b}_1 & \mathbf{a}_2^\top \mathbf{b}_2 & \cdots & \mathbf{a}_2^\top \mathbf{b}_p \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_m^\top \mathbf{b}_1 & \mathbf{a}_m^\top \mathbf{b}_2 & \cdots & \mathbf{a}_m^\top \mathbf{b}_p \end{bmatrix}\end{aligned}$$

Matrix–matrix products

Alternatively, we can represent \mathbf{A} by columns, and \mathbf{B} by rows, which leads to the interpretation of \mathbf{AB} as a sum of outer products.

Symbolically,

$$\mathbf{C} = \mathbf{AB} = \left[\begin{array}{c|c|c|c} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{array} \right] \left[\begin{array}{ccc} - & \mathbf{b}_1^\top & - \\ - & \mathbf{b}_2^\top & - \\ & \vdots & \\ - & \mathbf{b}_n^\top & - \end{array} \right] = \sum_{i=1}^n \mathbf{a}_i \mathbf{b}_i^\top$$

Put another way, \mathbf{AB} is equal to the sum, over all i , of the outer product of the i th column of \mathbf{A} and the i th row of \mathbf{B} . Since, in this case, $\mathbf{a}_i \in \mathbb{R}^m$ and $\mathbf{b}_i \in \mathbb{R}^p$, the dimension of the outer product $\mathbf{a}_i \mathbf{b}_i^\top$ is $m \times p$, which coincides with the dimension of \mathbf{C} .

Matrix–matrix products

We can also view matrix-matrix multiplication as a set of matrix-vector products. Specifically, if we represent \mathbf{B} by columns, we can view the columns of \mathbf{C} as matrix-vector products between \mathbf{A} and the columns of \mathbf{B} . Symbolically,

$$\mathbf{C} = \mathbf{AB} = \mathbf{A} \left[\begin{array}{c|c|c} | & | & \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & | & \end{array} \right] = \left[\begin{array}{c|c|c} | & | & \\ \mathbf{Ab}_1 & \mathbf{Ab}_2 & \cdots & \mathbf{Ab}_p \\ | & | & \end{array} \right].$$

Here the i th column of \mathbf{C} is given by the matrix-vector product with the vector on the right, $\mathbf{c}_i = \mathbf{Ab}_i$. These matrix-vector products can in turn be interpreted using both viewpoints given in the previous subsection.

Matrix–matrix products

Finally, we have the analogous viewpoint, where we represent \mathbf{A} by rows, and view the rows of \mathbf{C} as the matrix-vector product between the rows of \mathbf{A} and the matrix \mathbf{B} . Symbolically,

$$\mathbf{C} = \mathbf{AB} = \begin{bmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ & \vdots & \\ - & \mathbf{a}_m^\top & - \end{bmatrix} \mathbf{B} = \begin{bmatrix} - & \mathbf{a}_1^\top \mathbf{B} & - \\ - & \mathbf{a}_2^\top \mathbf{B} & - \\ & \vdots & \\ - & \mathbf{a}_m^\top \mathbf{B} & - \end{bmatrix}$$

It may seem like overkill to dissect matrix multiplication to such a large degree, especially when all these viewpoints follow immediately from the initial definition we gave (in about a line of math) at the beginning of this section. However, virtually all of linear algebra deals with matrix multiplications of some kind, and it is worthwhile to spend some time trying to develop an intuitive understanding of the viewpoints presented here.

1 Section 1: Introduction

- Notation
- Vector Space
- Norms
- Properties
- Special Types

2 Section 2: Matrix Multiplication

- Vector–vector products
- Matrix–vector products
- Matrix–matrix products
- **Manipulating data matrices**

Manipulating data matrices

Suppose \mathbf{X} is an $N \times D$ matrix. We can sum across the rows by premultiplying by a $1 \times N$ matrix of ones to create a $1 \times D$ matrix:

$$\mathbf{1}_N^\top \mathbf{X} = \left(\sum_n x_{n1} \quad \cdots \quad \sum_n x_{nD} \right)$$

Hence the mean of the data vectors is given by

$$\bar{\mathbf{x}}^\top = \frac{1}{N} \mathbf{1}_N^\top \mathbf{X}$$

We can sum across the columns by postmultiplying by a $D \times 1$ matrix of ones to create a $N \times 1$ matrix:

$$\mathbf{X} \mathbf{1}_D = \begin{pmatrix} \sum_d x_{1d} \\ \vdots \\ \sum_d x_{Nd} \end{pmatrix}$$

Manipulating data matrices

We can sum all entries in a matrix by pre and post multiplying by a vector of 1s:

$$\mathbf{1}_N^\top \mathbf{X} \mathbf{1}_D = \sum_{ij} X_{ij}$$

Hence the overall mean is given by

$$\bar{x} = \frac{1}{ND} \mathbf{1}_N^\top \mathbf{X} \mathbf{1}_D$$

We often want to scale rows or columns of a data matrix (e.g., to standardize them). We now show how to write this in matrix notation.

Manipulating data matrices

If we pre-multiply \mathbf{X} by a diagonal matrix $\mathbf{S} = \text{diag}(\mathbf{s})$, where \mathbf{s} is an N -vector, then we just scale each row of \mathbf{X} by the corresponding scale factor in \mathbf{s} :

$$\begin{aligned}\text{diag}(\mathbf{s})\mathbf{X} &= \begin{pmatrix} s_1 & \cdots & 0 \\ & \ddots & \\ 0 & \cdots & s_N \end{pmatrix} \begin{pmatrix} x_{1,1} & \cdots & x_{1,D} \\ & \ddots & \\ x_{N,1} & \cdots & x_{N,D} \end{pmatrix} = \\ &= \begin{pmatrix} s_1 x_{1,1} & \cdots & s_1 x_{1,D} \\ & \ddots & \\ s_N x_{N,1} & \cdots & s_N x_{N,D} \end{pmatrix}\end{aligned}$$

Manipulating data matrices

If we post-multiply \mathbf{X} by a diagonal matrix $\mathbf{S} = \text{diag}(\mathbf{s})$, where \mathbf{s} is a D -vector, then we just scale each column of \mathbf{X} by the corresponding element in \mathbf{s} .

$$\mathbf{X} \text{diag}(\mathbf{s}) = \begin{pmatrix} x_{1,1} & \cdots & x_{1,D} \\ & \ddots & \\ x_{N,1} & \cdots & x_{N,D} \end{pmatrix} \begin{pmatrix} s_1 & \cdots & 0 \\ & \ddots & \\ 0 & \cdots & s_D \end{pmatrix} =$$
$$\begin{pmatrix} s_1 x_{1,1} & \cdots & s_D x_{1,D} \\ & \ddots & \\ s_1 x_{N,1} & \cdots & s_D x_{N,D} \end{pmatrix}$$

Manipulating data matrices

Thus we can rewrite the standardization operation in matrix form as follows:

$$\text{standardize}(\mathbf{X}) = \left(\mathbf{X} - \mathbf{1}_N \boldsymbol{\mu}^T \right) \text{diag}(\boldsymbol{\sigma})^{-1}$$

where $\boldsymbol{\mu} = \bar{\mathbf{X}}$ is the empirical mean, and $\boldsymbol{\sigma}$ is a vector of the empirical standard deviations.

Manipulating data matrices

Sum Square Matrix

The sum of squares matrix is $D \times D$ matrix defined by

$$\mathbf{S}_0 \triangleq \mathbf{X}^\top \mathbf{X} = \sum_{n=1}^{N_D} \mathbf{x}_n \mathbf{x}_n^\top = \sum_{n=1}^{N_D} \begin{pmatrix} x_{n,1}^2 & \cdots & x_{n,1}x_{n,D} \\ & \ddots & \\ x_{n,D}x_{n,1} & \cdots & x_{n,D}^2 \end{pmatrix}$$

The scatter matrix is a $D \times D$ matrix defined by

$$\mathbf{S}_{\bar{\mathbf{x}}} \triangleq \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^\top = \left(\sum_n \mathbf{x}_n \mathbf{x}_n^\top \right) - N \bar{\mathbf{x}} \bar{\mathbf{x}}^\top$$

Manipulating data matrices

We see that this is the sum of squares matrix applied to the mean-centered data. More precisely, define $\tilde{\mathbf{X}}$ to be a version of \mathbf{X} where we subtract the mean $\bar{\mathbf{x}} = \frac{1}{N}\mathbf{X}^\top \mathbf{1}_N$ off every row. Hence we can compute the centered data matrix using

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top = \mathbf{X} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{X} = \mathbf{C}_N \mathbf{X}$$

where

$$\mathbf{C}_N \triangleq \mathbf{I}_N - \frac{1}{N} \mathbf{J}_N$$

is the centering matrix, and $\mathbf{J}_N = \mathbf{1}_N \mathbf{1}_N^\top$ is a matrix of all 1s. The scatter matrix can now be computed as follows:

$$\mathbf{S}_{\bar{\mathbf{x}}} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \mathbf{X}^\top \mathbf{C}_N^\top \mathbf{C}_N \mathbf{X} = \mathbf{X}^\top \mathbf{C}_N \mathbf{X}$$

Manipulating data matrices

Gram Matrix

The $N \times N$ matrix \mathbf{XX}^\top is a matrix of inner products called the Gram matrix:

$$\mathbf{K} \triangleq \mathbf{XX}^\top = \begin{pmatrix} \mathbf{x}_1^\top \mathbf{x}_1 & \cdots & \mathbf{x}_1^\top \mathbf{x}_N \\ & \ddots & \\ \mathbf{x}_n^\top \mathbf{x}_1 & \cdots & \mathbf{x}_N^\top \mathbf{x}_N \end{pmatrix}$$

Sometimes we want to compute the inner products of the mean-centered data vectors, $\tilde{\mathbf{K}} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$. However, if we are working with a feature similarity matrix instead of raw features, we will only have access to \mathbf{K} , not \mathbf{X} . Fortunately, we can compute $\tilde{\mathbf{K}}$ from \mathbf{K} using the double centering trick:

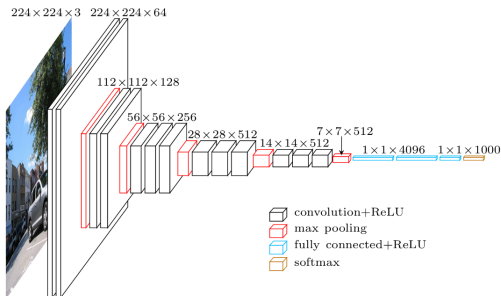
$$\tilde{\mathbf{K}} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top = \mathbf{C}_N \mathbf{K} \mathbf{C}_N = \mathbf{K} - \frac{1}{N} \mathbf{J} \mathbf{K} - \frac{1}{N} \mathbf{K} \mathbf{J} + \frac{1}{N^2} \mathbf{J} \mathbf{K} \mathbf{J}$$

Manipulating data matrices

This subtracts the row means and column means from \mathbf{K} , and adds back the global mean that gets subtracted twice, so that both row means and column means of $\tilde{\mathbf{K}}$ are equal to zero.

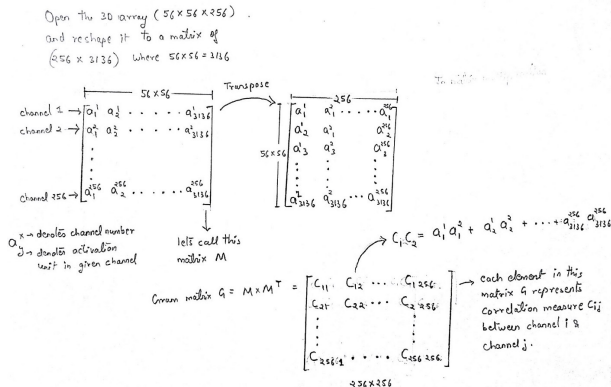
Manipulating data matrices

Gram matrix is often used in Deep learning, such as NN Transfer Style. The principle is simple: we define two distances, one for the content (DC) and one for the style (DS). DC measures how different the content is between two images while DS measures how different the style is between two images. Then, we take a third image, the input, and transform it to minimize both its content-distance with the content-image and its style-distance with the style-image. Actually we use Gram matrix to measure the DS part.



Manipulating data matrices

Calculation Process:



Manipulating data matrices

Distance Matrix

Let \mathbf{X} be $N_x \times D$ datamatrix, and \mathbf{Y} be another $N_y \times D$ datamatrix. We can compute the squared pairwise distances between these using

$$\mathbf{D}_{ij} = (\mathbf{x}_i - \mathbf{y}_j)^\top (\mathbf{x}_i - \mathbf{y}_j) = \|\mathbf{x}_i\|^2 - 2\mathbf{x}_i^\top \mathbf{y}_j + \|\mathbf{y}_j\|^2$$

Let us now write this in matrix form. Let

$\hat{\mathbf{x}} = [\|\mathbf{x}_1\|^2; \dots; \|\mathbf{x}_{N_x}\|^2] = \text{diag}(\mathbf{X}\mathbf{X}^\top)$ be a vector where each element is the squared norm of the examples in \mathbf{X} , and define $\hat{\mathbf{y}}$ similarly. Then we have

$$\mathbf{D} = \hat{\mathbf{x}}\mathbf{1}_{N_y}^\top - 2\mathbf{X}\mathbf{Y}^\top + \mathbf{1}_{N_x}\hat{\mathbf{y}}^\top$$

In the case that $\mathbf{X} = \mathbf{Y}$, we have

$$\mathbf{D} = \hat{\mathbf{x}}\mathbf{1}_N^\top - 2\mathbf{X}\mathbf{X}^\top + \mathbf{1}_N\hat{\mathbf{x}}^\top$$

Einstein Summation

Einstein summation, or einsum for short, is a notational shortcut for working with tensors. For example, instead of writing matrix multiplication as $C_{ij} = \sum_k A_{ik} B_{kj}$, we can just write it as $C_{ij} = A_{ik} B_{kj}$, where we drop the \sum_k .