

Linear Algebra Background

Winter 2023

Clayton Scott

1 Overview

These notes review concepts that are critical for machine learning. Many of the concepts discussed below can be defined in more general settings, but for concreteness I focus on *Euclidean* vectors. I assume the reader is familiar with basic matrix and vector operations, as well as basic set notation.

2 Dot Product, Norm, and Outer Product

Let

$$\mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_d \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix}$$

be two vectors. The *dot product* of \mathbf{u} and \mathbf{v} is

$$\begin{aligned} \langle \mathbf{u}, \mathbf{v} \rangle &= u_1 v_1 + \cdots + u_d v_d \\ &= \sum_{i=1}^d u_i v_i. \end{aligned}$$

The dot product is an example of an inner product. We'll talk about the more general definition of an inner product later in the course.

The (*Euclidean*) *norm* of a vector $\mathbf{u} \in \mathbb{R}^n$ is defined by

$$\|\mathbf{u}\| := \sqrt{\langle \mathbf{u}, \mathbf{u} \rangle} = \left(\sum_{j=1}^n u_j^2 \right)^{\frac{1}{2}}$$

Example 1. If $\mathbf{u} = \begin{bmatrix} 4 \\ 2 \\ -7 \end{bmatrix}$, then $\|\mathbf{u}\| = \sqrt{16 + 4 + 49} = \sqrt{69} \approx 8.3$.

The Cauchy-Schwartz inequality is a key property of inner products:

Theorem 1. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$, with equality iff \mathbf{x} and \mathbf{y} are scalar multiples of one another (i.e., they are linearly dependent).

We say \mathbf{u} and \mathbf{v} are *orthogonal* if $\mathbf{u} \neq \mathbf{0} \neq \mathbf{v}$ (where $\mathbf{0}$ denotes the zero vector), and $\langle \mathbf{u}, \mathbf{v} \rangle = 0$. See Figure 1.

The dot product of \mathbf{u} and \mathbf{v} can be expressed using matrix multiplication as

$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{v}.$$

The product $\mathbf{u}\mathbf{v}^T$ is a $d \times d$ matrix called the *outer product* of \mathbf{u} and \mathbf{v} .

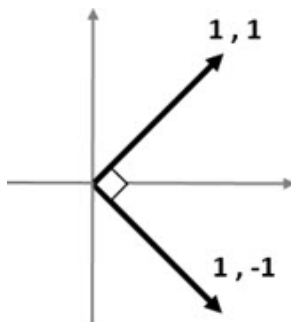


Figure 1: Two orthogonal vectors.

3 Matrix-Vector Multiplication

Let

$$\mathbf{A} = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \in \mathbb{R}^n.$$

By the definition of matrix multiplication,

$$\mathbf{Ab} = \begin{bmatrix} \langle \mathbf{A}_{1,:}, \mathbf{b} \rangle \\ \vdots \\ \langle \mathbf{A}_{m,:}, \mathbf{b} \rangle \end{bmatrix} \in \mathbb{R}^m,$$

where $\mathbf{A}_{i,:}$ denotes the i^{th} row of \mathbf{A} .

However, there is another (very useful) way to represent \mathbf{Ab} . Define

$$\mathbf{a}_j = \mathbf{A}_{:,j},$$

the j^{th} column of \mathbf{A} . Then

$$\mathbf{Ab} = \sum_{j=1}^n \underbrace{b_j}_{\text{scalar}} \underbrace{\mathbf{a}_j}_{\text{vector}}.$$

To see this, it helps to write

$$\begin{aligned} \sum_j b_j \mathbf{a}_j &= b_1 \begin{bmatrix} a_{11} \\ \vdots \\ a_{m1} \end{bmatrix} + \cdots + b_n \begin{bmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{bmatrix} \\ &= \begin{bmatrix} \sum b_j a_{1j} \\ \vdots \\ \sum b_j a_{mj} \end{bmatrix} \\ &= \mathbf{Ab}. \end{aligned}$$

4 Span

Now consider an arbitrary collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^m$. A *linear combination* of these vectors is any vector of the form

$$\sum_{j=1}^n x_j \mathbf{a}_j,$$

where $x_1, \dots, x_n \in \mathbb{R}$. The *span* of $\mathbf{a}_1, \dots, \mathbf{a}_n$ is the set of all linear combinations of $\mathbf{a}_1, \dots, \mathbf{a}_n$. If we denote by \mathbf{A} the matrix whose j^{th} column is \mathbf{a}_j , i.e.,

$$\mathbf{A} = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \cdots & \mathbf{a}_n \\ | & & | \end{bmatrix},$$

and we use the insight of the previous section, the (*linear*) *span* of $\mathbf{a}_1, \dots, \mathbf{a}_n$ is

$$\text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_n\} = \{\mathbf{y} = \mathbf{A}\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\}.$$

This collection of vectors is also referred to as the *column span* or *image* of \mathbf{A} , and denoted $\text{colspan}(\mathbf{A})$. See Figure 2.

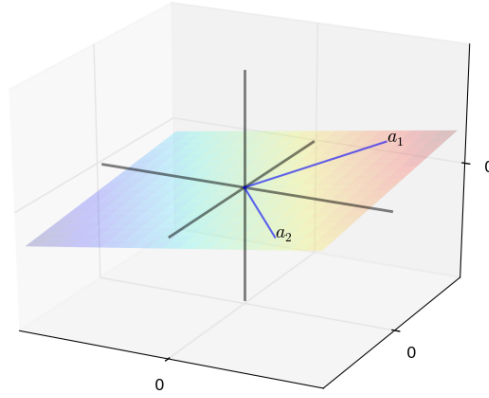


Figure 2: Linear span of two vectors.

5 Linear Independence

Vectors $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^m$ are *linearly independent* iff the following implication holds:

$$\sum_{j=1}^n x_j \mathbf{a}_j = \mathbf{0} \implies x_j = 0 \forall j.$$

We can also think of linear independence in terms of the matrix \mathbf{A} as in the previous section. In particular, for any matrix \mathbf{A} , the *nullspace* of \mathbf{A} is the set

$$N(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^m | \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

If

$$\mathbf{A} = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \cdots & \mathbf{a}_n \\ | & & | \end{bmatrix},$$

then $\mathbf{a}_1, \dots, \mathbf{a}_n$ are linearly independent iff $N(\mathbf{A}) = \{\mathbf{0}\}$. This is simply a restatement of the definition.

If a set of vectors is not linearly independent, the vectors are said to be *linearly dependent*.

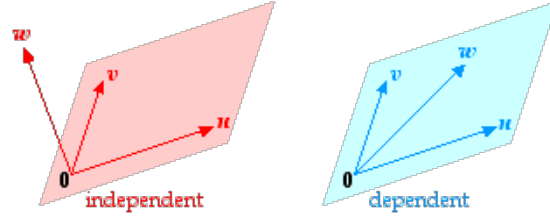


Figure 3: Left: Linearly independent vectors. Right: Linearly dependent vectors.

6 Subspaces

A *subspace* is a set of vectors that is closed under scalar multiplication and vector addition. In particular, a set $S \subseteq \mathbb{R}^m$ is a subspace provided the following two properties hold:

1. For all $\alpha \in \mathbb{R}$ and $\mathbf{u} \in S$, $\alpha\mathbf{u} \in S$.
2. For all $\mathbf{u}, \mathbf{v} \in S$, $\mathbf{u} + \mathbf{v} \in S$.

Example 2. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. Then $\text{colspan}(\mathbf{A})$ is a subspace of \mathbb{R}^m , and $N(\mathbf{A})$ is a subspace of \mathbb{R}^n .

A *basis* for a subspace $S \subseteq \mathbb{R}^d$ is a finite collection of vectors $B \subseteq S$ satisfying any one of the following equivalent properties:

1. B is a minimal spanning set (i.e., a spanning set that is no longer a spanning set if some element is removed)
2. B is a maximal linearly independent set (i.e., a linearly independent set that is no longer a linearly independent set if any vector is added)
3. $\text{span}(B) = S$ and B is a linearly independent set
4. Every element of S can be written as a unique linear combination of elements of B .

Example 3. Consider the matrix

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

The vectors $\mathbf{a}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ and $\mathbf{a}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ are a basis of the image of \mathbf{A} . To see this, we can use the 3rd characterization of bases: The vectors are clearly linearly independent. Furthermore, they span $\text{colspan}(\mathbf{A})$. To see this, observe that $\mathbf{a}_3 = \mathbf{a}_1 + \mathbf{a}_2$. Then for any \mathbf{x} ,

$$\begin{aligned} \mathbf{A}\mathbf{x} &= x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + x_3\mathbf{a}_3 \\ &= (x_1 + x_3)\mathbf{a}_1 + (x_2 + x_3)\mathbf{a}_2. \\ &\in \text{span}\{\mathbf{a}_1, \mathbf{a}_2\}. \end{aligned}$$

By the third characterization of bases, these two vectors constitute a basis.

The *dimension* of a subspace is the cardinality of a basis. The set $\{\mathbf{0}\}$ is a subspace, and its dimension is taken to be zero.

Example 4. $\dim(\mathbb{R}^n) = n$, with the *standard basis*

$$\mathbf{e}_i = [0, \dots, 0, \underbrace{1}_{i^{\text{th}} \text{ position}}, 0, \dots, 0]^{\top}, \quad i = 1, \dots, n$$

as one possible basis.

7 Rank and Nullity

Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, define the *rank* of \mathbf{A}

$$\text{rank}(\mathbf{A}) = \dim(\text{colspan}(\mathbf{A})),$$

and the *nullity* of \mathbf{A}

$$\text{nullity}(\mathbf{A}) = \dim(N(\mathbf{A})).$$

If $\text{rank}(\mathbf{A}) = \min(m, n)$, \mathbf{A} is said to have *full rank*. If $n < m$, then \mathbf{A} has full rank iff the columns of \mathbf{A} are linearly independent. In this case, we say \mathbf{A} has *full column rank*.

Rank and nullity are related by the following result:

Theorem 2 (Rank-Nullity Theorem). *If $\mathbf{A} \in \mathbb{R}^{m \times n}$, then $\text{rank}(\mathbf{A}) + \text{nullity}(\mathbf{A}) = n$.*

Example 5. In the example of the previous section, $\text{rank}(\mathbf{A}) = 2$ and therefore $\text{nullity}(\mathbf{A}) = 1$.

8 Inverses

A $d \times d$ square matrix \mathbf{A} is invertible when there exists another square matrix \mathbf{B} such that $\mathbf{AB} = \mathbf{BA} = \mathbf{I}$, the identity matrix. In this case the matrix \mathbf{B} is denoted \mathbf{A}^{-1} . The matrix \mathbf{A} is invertible iff any one of the following equivalent conditions hold:

1. $\text{rank}(\mathbf{A}) = d$
2. The columns of \mathbf{A} are linearly independent
3. $\text{nullity}(\mathbf{A}) = 0$
4. If $\mathbf{Ax} = \mathbf{0}$, then $\mathbf{x} = \mathbf{0}$.
5. $\det(\mathbf{A}) \neq 0$
6. For all $\mathbf{b} \in \mathbb{R}^d$, the equation $\mathbf{Ax} = \mathbf{b}$ has a unique solution.

9 Orthogonal and Orthonormal Sets

The nonzero vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ comprise an *orthogonal set* if $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$ for all $i \neq j$. If in addition $\|\mathbf{u}_i\| = 1$ for all i , we say $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ is an *orthonormal set*.

Example 6.

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

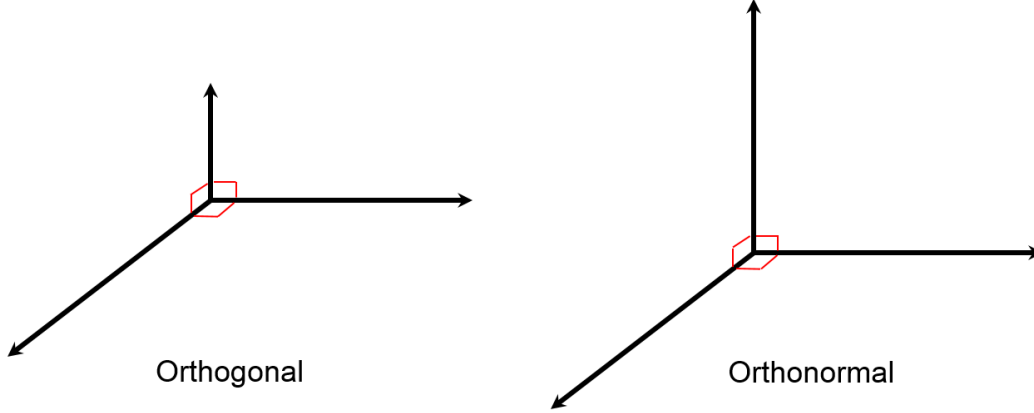
are orthonormal.

Every orthogonal set is linearly independent. To see this, suppose $\mathbf{u}_1, \dots, \mathbf{u}_n$ are orthogonal and $c_1\mathbf{u}_1 + \dots + c_n\mathbf{u}_n = \mathbf{0}$ for some scalars $c_1, \dots, c_n \in \mathbb{R}$. Then for each i

$$\begin{aligned} 0 &= \langle \mathbf{0}, \mathbf{u}_i \rangle = \langle c_1\mathbf{u}_1 + \dots + c_n\mathbf{u}_n, \mathbf{u}_i \rangle \\ &= c_i \|\mathbf{u}_i\|^2 \end{aligned}$$

which implies $c_i = 0$ since $\|\mathbf{u}_i\| > 0$.

If $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^n$ are an orthogonal set, they are a basis of \mathbb{R}^n , by the second characterization of basis: they are LI and clearly maximal since $\dim(\mathbb{R}^n) = n$.



If $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ is an orthonormal set and a basis of a subspace S , it is referred to as an *orthonormal basis* (ONB) of that subspace.

An *orthogonal matrix* is a square matrix \mathbf{U} such that $\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$. That is, the transpose of \mathbf{U} is its inverse. Note that $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ implies the columns of \mathbf{U} are an ONB, and $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ implies the rows of \mathbf{U} are also an ONB.

Example 7. The rotation matrix $\begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$ is an orthogonal matrix.

10 Eigenvalues and Eigenvectors

Let $\mathbf{A} \in \mathbb{R}^{d \times d}$. If

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

for some $\lambda \in \mathbb{R}$ and $\mathbf{u} \in \mathbb{R}^d$, we say λ is an *eigenvalue* of \mathbf{A} and \mathbf{u} is a corresponding *eigenvector*. If \mathbf{A} is *symmetric*, i.e., $\mathbf{A} = \mathbf{A}^T$, we can characterize \mathbf{A} in terms of its eigenvalues and eigenvectors.

Theorem 3 (Spectral Theorem). *If $\mathbf{A} \in \mathbb{R}^{d \times d}$ is symmetric, then*

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

where \mathbf{U} is an orthogonal matrix with real entries and $\mathbf{\Lambda}$ is a diagonal matrix with real entries.

What is the connection to eigenvalues and eigenvectors? Multiplying $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ on the right by \mathbf{U} , we have

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}.$$

Let

$$\mathbf{U} = \begin{bmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_d \\ | & & | \end{bmatrix}, \quad \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_d \end{bmatrix}.$$

If we look at the matrix equation $\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}$ one column at a time, we have

$$\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i, i = 1, \dots, n.$$

Thus, λ_i are eigenvalues of \mathbf{A} , and \mathbf{u}_i are corresponding eigenvectors. Since \mathbf{U} is an orthogonal matrix, $\mathbf{u}_1, \dots, \mathbf{u}_d$ are an orthonormal basis of \mathbb{R}^d . Thus, the spectral theorem implies the existence of an ONB consisting of eigenvectors of \mathbf{A} .

An important identity is

$$\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T.$$

This can be verified just by calculating entries on both sides and checking that they agree. For example, the $(1, 1)$ entry is $\sum_{i=1}^d \lambda_i u_{i1}^2$. This identity expresses \mathbf{A} as a sum of rank 1 outer products. $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ and $\sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ are referred to as the *spectral* or *eigenvalue decomposition* of \mathbf{A} . $\lambda_1, \dots, \lambda_d$ constitute the *spectrum* of \mathbf{A} .

The determinant of a matrix is the product of its eigenvalues, while the trace of a matrix (the sum of the diagonal entries) is the sum of its eigenvalues.

11 Positive (Semi-)Definite Matrices

Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a square matrix. We say

- \mathbf{A} is *positive semi-definite* if $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^d$
- \mathbf{A} is *positive definite* if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$.

Clearly if \mathbf{A} is PD, it is also PSD.

PD and PSD matrices arise frequently in machine learning. For example,

- Gram matrices, which have the form

$$\begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_n \rangle \\ \vdots & \ddots & \vdots \\ \langle \mathbf{x}_n, \mathbf{x}_1 \rangle & \cdots & \langle \mathbf{x}_n, \mathbf{x}_n \rangle \end{bmatrix}$$

for some vectors $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^d$

- Covariance matrices
- Kernel matrices
- Hessian matrices of convex functions.

All of the above example are symmetric, and all of the PD/PSD matrices we will encounter will be symmetric. By the spectral theorem, such matrices have an eigenvalue decomposition with real eigenvalues. The following properties characterize PSD and PD matrices in terms of their spectrum.

- \mathbf{A} is PSD iff $\lambda_i \geq 0 \ \forall i$.
- \mathbf{A} is PD iff $\lambda_i > 0 \ \forall i$.

12 Orthogonal Complements

Let $S \subseteq \mathbb{R}^d$ be a subspace. The set

$$S^\perp := \{\mathbf{v} \mid \langle \mathbf{u}, \mathbf{v} \rangle = 0 \text{ for all } \mathbf{u} \in S\}$$

is called the orthogonal complement of S . It can easily be shown to be a subspace itself. See Figure 4.

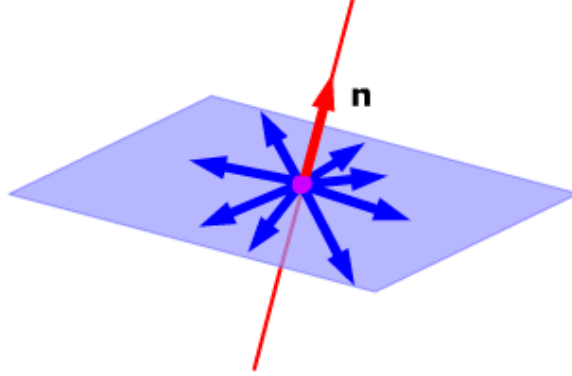


Figure 4: Illustration of orthogonal complement

13 Projections

Let $\mathbf{a}_1, \dots, \mathbf{a}_k \in \mathbb{R}^d$ be linearly independent column vectors. Denote

$$\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_k] \quad (d \times k).$$

The linear span of $\mathbf{a}_1, \dots, \mathbf{a}_k$ is the column span of \mathbf{A} , written $\langle \mathbf{A} \rangle$ for brevity.

The *projection* onto the subspace $S = \langle \mathbf{A} \rangle$ is the mapping $\Pi_{\mathbf{A}}$ that sends $\mathbf{x} \in \mathbb{R}^d$ to the closest point in $\langle \mathbf{A} \rangle$. Every point in $\langle \mathbf{A} \rangle$ equals $\mathbf{A}\boldsymbol{\theta}$ for some $\boldsymbol{\theta} \in \mathbb{R}^k$. Therefore, $\Pi_{\mathbf{A}}\mathbf{x} = \mathbf{A}\hat{\boldsymbol{\theta}}$, where

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^k} \|\mathbf{x} - \mathbf{A}\boldsymbol{\theta}\|^2.$$

The solution is

$$\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}$$

where $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ is called the pseudoinverse of \mathbf{A} . Since \mathbf{A} has full rank, $\mathbf{A}^T \mathbf{A}$ is invertible (exercise). Therefore,

$$\Pi_{\mathbf{A}} \mathbf{x} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}.$$

Since \mathbf{x} was arbitrary, we conclude that $\Pi_{\mathbf{A}}$ is the $d \times d$ matrix

$$\Pi_{\mathbf{A}} := \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T,$$

which is called the *projection matrix* for \mathbf{A} . Projection matrices are symmetric and positive-semidefinite.

If $\mathbf{a}_1, \dots, \mathbf{a}_k$ are orthonormal, i.e.,

$$\langle \mathbf{a}_i, \mathbf{a}_j \rangle = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases},$$

then

$$\Pi_{\mathbf{A}} = \mathbf{A} \mathbf{A}^T \quad (d \times d).$$

Since the projection depends only on the subspace S and not the matrix used to represent it, we can assume the columns of \mathbf{A} are orthonormal without loss of generality.

The *orthogonality principle* states that $\forall \mathbf{x}, \mathbf{x} - \Pi_{\mathbf{A}} \mathbf{x}$ is orthogonal to every element of $\langle \mathbf{A} \rangle$. That is, $\mathbf{x} - \Pi_{\mathbf{A}} \mathbf{x} \in \langle \mathbf{A} \rangle^\perp$ (see Figure 5 (b)). To see this, let $\mathbf{A}\boldsymbol{\theta}$ denote an arbitrary element of $\langle \mathbf{A} \rangle$. Then

$$\begin{aligned} \langle \mathbf{A}\boldsymbol{\theta}, \mathbf{x} - \Pi_{\mathbf{A}} \mathbf{x} \rangle &= \boldsymbol{\theta}^T \mathbf{A}^T (\mathbf{x} - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}) \\ &= \boldsymbol{\theta}^T (\mathbf{A}^T \mathbf{x} - \mathbf{A}^T \mathbf{x}) \\ &= 0. \end{aligned}$$

The *projection theorem* states that for any subspace $S \subseteq \mathbb{R}^d$, $\mathbb{R}^d = S \oplus S^\perp$. In other words, every $\mathbf{x} \in \mathbb{R}^d$ can be written uniquely as $\mathbf{u} + \mathbf{v}$ where $\mathbf{u} \in S$ and $\mathbf{v} \in S^\perp$. Existence of \mathbf{u} and \mathbf{v} follows from the orthogonality principle by taking $\mathbf{u} = \Pi_A \mathbf{x}$ (where A is such that $S = \langle A \rangle$) and $\mathbf{v} = \mathbf{x} - \mathbf{u}$. To see uniqueness, suppose $\mathbf{x} = \mathbf{u}' + \mathbf{v}'$. Then $\mathbf{u} - \mathbf{u}' = \mathbf{v}' - \mathbf{v}$. But the only common element of S and S^\perp is the zero vector, so $\mathbf{u} = \mathbf{u}'$ and $\mathbf{v} = \mathbf{v}'$.

Projection matrices are *idempotent*, which means $\Pi_A^2 = \Pi_A$. This follows from

$$\begin{aligned}\Pi_A^2 &= \Pi_A \cdot \Pi_A \\ &= A(A^T A)^{-1} A^T \cdot A(A^T A)^{-1} A^T \\ &= A(A^T A)^{-1} A^T \\ &= \Pi_A.\end{aligned}$$

Intuitively, the second projection has no effect because $\Pi_A \mathbf{x} \in \langle A \rangle$ already.

Remark 1. In the above treatment of projections, we defined projections in terms of the closest point property, from which we derived the orthogonality condition and the projection theorem. Projections can be defined in a much more general setting, namely projections onto closed subspaces of Hilbert spaces. In this more general setting, one first proves the projection theorem and uses it to define projections, and then deduces the orthogonality principle and closest point property as corollaries. Projection in Hilbert space is a very important concept in applied mathematics and data analysis.

Exercises

1. (★) Prove the Pythagorean theorem: If $\langle \mathbf{u}, \mathbf{v} \rangle = 0$, then $\|\mathbf{u} + \mathbf{v}\|^2 = \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2$.
2. (★) Verify that the column span and nullspace of a matrix are subspaces.
3. (★) For the matrix in Example 3, find a basis for the nullspace.
4. (★) Show that if \mathbf{u} and \mathbf{v} are two vectors, then the outer product $\mathbf{u}\mathbf{v}^T$ has rank one.
5. (★★) Show that the sum of k rank-one matrices has rank at most k .
6. (★★) Let $k < d$. Show that if $A \in \mathbb{R}^{d \times k}$ is full rank, then $A^T A$ is invertible.
7. (★★) Use the projection theorem to prove the rank plus nullity theorem. *Hint:* The row-rank of a matrix is equal to the column rank, which you may assume.
8. (★) Show that if λ is an eigenvalue of a matrix, then the set of all associated eigenvectors constitute a subspace. This subspace is called the *eigenspace* of λ and its dimension is called the (geometric) multiplicity of λ .
9. (a) (★) Show that if U is an orthogonal matrix, then for all $\mathbf{x} \in \mathbb{R}^d$, $\|\mathbf{x}\| = \|U\mathbf{x}\|$, where the norm is the Euclidean norm.
(b) (★★) Show that all 2×2 orthogonal matrices have the form

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix}$$

for some θ . Give a geometric interpretation of the effect of these two transformations.

- (c) (★★) An ellipse in \mathbb{R}^2 can be expressed in the form

$$\left\{ \mathbf{x} \mid (\mathbf{x} - \mathbf{c})^T \mathbf{A} (\mathbf{x} - \mathbf{c}) = r^2 \right\},$$

where $\mathbf{c} \in \mathbb{R}^2$, $r > 0$, and $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ where \mathbf{U} is orthogonal and $\mathbf{\Lambda}$ is diagonal and positive definite. Choose $\mathbf{c}, r, \mathbf{U}$, and $\mathbf{\Lambda}$ such that the major axis has length 3, the minor axis has length 1, and the center of the ellipse is at the point $[3 \ -1]^T$, and the major axis makes an angle of $+\pi/6$ radians with the positive x -axis.

10. (★★) Let \mathbf{A} be a symmetric matrix. Show that \mathbf{A} is invertible iff all of its eigenvalues are nonzero, and express the spectral decomposition of \mathbf{A}^{-1} in terms of the spectral decomposition of \mathbf{A} . Conclude that if \mathbf{A} is PD, then so is \mathbf{A}^{-1} .
11. Show that the following types of matrices are PSD:
 - (a) (★) Any matrix of the form $\mathbf{A}^T \mathbf{A}$, where \mathbf{A} is an arbitrary matrix.
 - (b) (★) Covariance matrices, i.e., matrices of the form $\mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T]$, where \mathbf{X} is a random column vector.
 - (c) (★★) Gram matrices.
12. (☆☆) Show that a Gram matrix is PD iff the associated vectors are linearly independent.
13. (☆☆) Prove the Cauchy-Schwartz inequality by considering the determinant of the 2×2 Gram matrix associated to \mathbf{x} and \mathbf{y} .
14. (☆) Show that the matrix

$$\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

is positive definite but not symmetric.

15. (★★) Describe the eigenvalue decomposition of the projection matrix for a subspace S . What are the eigenvalues, what are their multiplicities, and what are the associated eigenspaces? (See an earlier problem for definitions of multiplicity and eigenspace).
16. (★★) Show that $(S^\perp)^\perp = S$, and use this fact together with the projection theorem to show that the projection onto S^\perp is given by $\mathbf{I} - \mathbf{\Pi}_A$, where $S = \langle \mathbf{A} \rangle$.

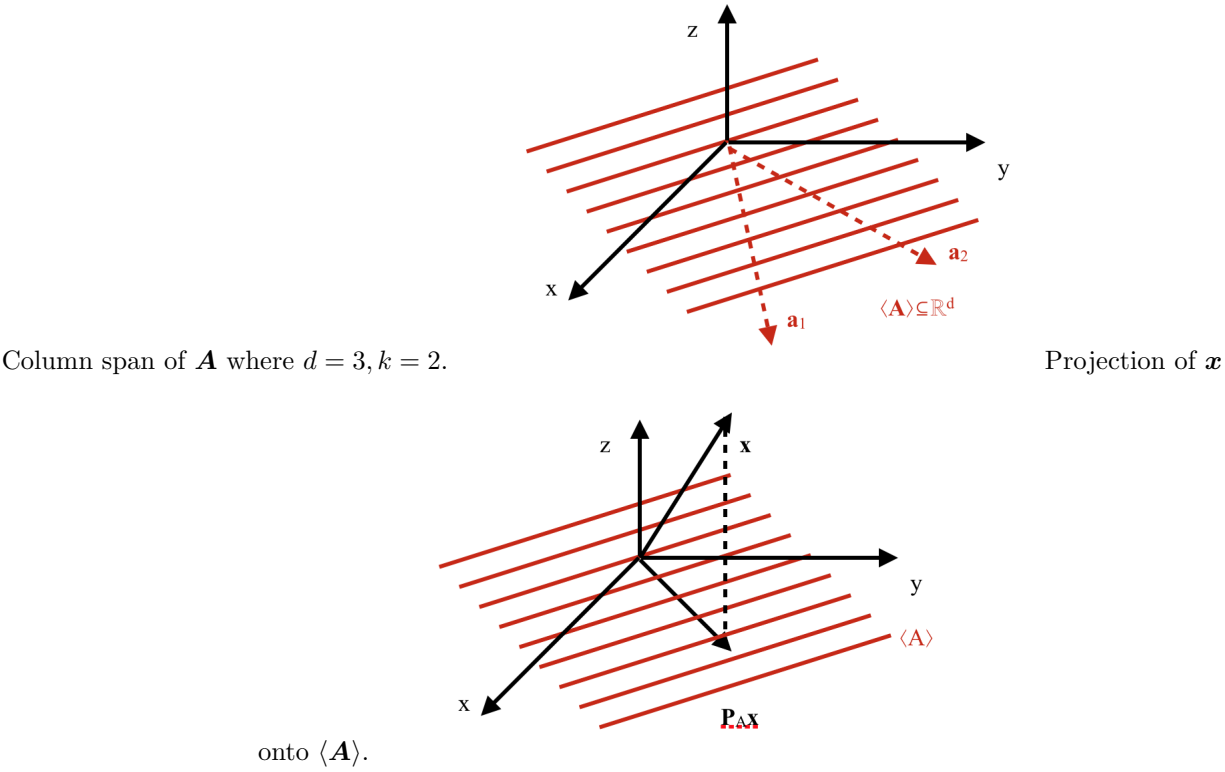


Figure 5: Linear span and projection