

Mathematics for Machine Learning

Prof Willie Brink

Applied Mathematics, Stellenbosch University

Lecture 1: Linear Algebra

Module information

Lecturer: Prof Willie Brink (wbrink@sun.ac.za)

Classes: Mondays and Wednesdays 9:30 to 12:30 in A403A

Textbook: <https://mml-book.com>

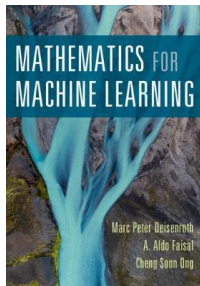
Assessment:

Week 1: assignment 1 (15%), quiz 1 (15%)

Week 2: assignment 2 (15%), quiz 2 (15%)

Week 3: assignment 3 (15%), final test (25%)

SUNLearn: <https://learn.sun.ac.za>



Contents of the module

Chapter 2: Linear Algebra

Chapter 03: Analytic Geometry

Chapter 04: Matrix Decompositions

Chapter 05: Vector Calculus

Chapter 06: Probability and Distributions

Chapter 07: Continuous Optimisation

Chapter 08: When Models Meet Data

Chapter 09: Linear Regression

Chapter 10: Dimensionality Reduction with Principal Component Analysis

Chapter 11: Density Estimation with Gaussian Mixture Models

Chapter 12: Classification with Support Vector Machines

2.1 Systems of linear equations

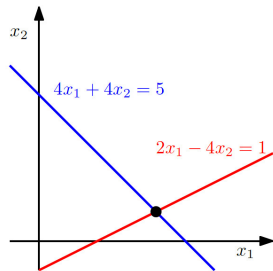
$$\begin{array}{rcccccccc} a_{1,1} & x_1 & + & a_{1,2} & x_2 & + & \dots & + & a_{1,n} & x_n & = & b_1 \\ a_{2,1} & x_1 & + & a_{2,2} & x_2 & + & \dots & + & a_{2,n} & x_n & = & b_2 \\ & \vdots & & & & & & & & & & \vdots \\ a_{m,1} & x_1 & + & a_{m,2} & x_2 & + & \dots & + & a_{m,n} & x_n & = & b_m \end{array}$$

$$\begin{bmatrix} a_{1,1} \\ a_{2,1} \\ \vdots \\ a_{m,1} \end{bmatrix} x_1 + \begin{bmatrix} a_{1,2} \\ a_{2,2} \\ \vdots \\ a_{m,2} \end{bmatrix} x_2 + \dots + \begin{bmatrix} a_{1,n} \\ a_{2,n} \\ \vdots \\ a_{m,n} \end{bmatrix} x_n = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

$$\begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \implies \mathbf{Ax} = \mathbf{b}$$

m equations, n unknowns

Can have no solution,
or exactly one solution,
or infinitely many.



2.2 Matrices

$\mathbb{R}^{m \times n}$ is the set of all real-valued matrices with m rows and n columns.

The **sum** of matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is computed elementwise.

The **product** of $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times k}$ is a matrix $\mathbf{C} \in \mathbb{R}^{m \times k}$ with $c_{i,j} = \sum_{\ell=1}^n a_{i,\ell} b_{\ell,j}$.

The Hadamard product of matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is computed elementwise.

Matrix multiplication is associative and distributive, but in general not commutative ($\mathbf{AB} \neq \mathbf{BA}$).

With \mathbf{I}_n the **identity matrix** in $\mathbb{R}^{n \times n}$, we have $\mathbf{I}_m \mathbf{A} = \mathbf{A}$ and $\mathbf{A} \mathbf{I}_n = \mathbf{A}$, $\forall \mathbf{A} \in \mathbb{R}^{m \times n}$.

The **inverse** of square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a matrix $\mathbf{B} = \mathbf{A}^{-1}$ such that $\mathbf{AB} = \mathbf{BA} = \mathbf{I}_n$.

- if the inverse exists, \mathbf{A} is called invertible / nonsingular (or regular)
- if the inverse doesn't exist, \mathbf{A} is called noninvertible / singular

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1} \quad \text{and, in general, } (\mathbf{A} + \mathbf{B})^{-1} \neq \mathbf{A}^{-1} + \mathbf{B}^{-1}.$$

The **transpose** of $\mathbf{A} \in \mathbb{R}^{m \times n}$ is $\mathbf{B} = \mathbf{A}^T \in \mathbb{R}^{n \times m}$, with $b_{i,j} = a_{j,i}$.

$$(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T \quad \text{and} \quad (\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T.$$

A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is **symmetric** if $\mathbf{A} = \mathbf{A}^T$.

If \mathbf{A} is invertible, then so is \mathbf{A}^T , and $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T = \mathbf{A}^{-T}$.

Scalar multiplication ($\lambda\mathbf{A}$) is calculated elementwise, and is associative and distributive.

2.3 Solving systems of linear equations

General approach to solve $\mathbf{Ax} = \mathbf{b}$:

1. find a particular solution to $\mathbf{Ax} = \mathbf{b}$
2. find all solutions to $\mathbf{Ax} = \mathbf{0}$
3. combine the solutions from steps 1 and 2

Gaussian elimination

Use **elementary transformations** that do not change the solution (row exchange, multiplying a row with nonzero constant, adding rows) to find a **row-echelon form**.

pivot: first nonzero element in a row from the left

staircase structure: every pivot is strictly to the right of the pivot above it

The **reduced row echelon form**, where every pivot is 1 and is the only nonzero entry in its column, eases steps 1 and 2 above.

1. Finding a **particular solution** to $\mathbf{Ax} = \mathbf{b}$:

Write $[\mathbf{A} \mid \mathbf{b}]$ in reduced row-echelon form (RREF).

Set free variables (not corresponding to pivots) to zero.

Easily solve for the basic variables (corresponding to pivots).

2. Finding a **general solution** to $\mathbf{Ax} = \mathbf{0}$:

Augment the RREF of \mathbf{A} with rows of the form $[0 \cdots 0 \ -1 \ 0 \cdots 0]$ so that we have 1 or -1 on the diagonal.

General solution: all linear combinations of the columns with -1 on the diagonal.

3. A general solution to $\mathbf{Ax} = \mathbf{b}$ will be the sum of steps 1 and 2.

Calculating the inverse of $\mathbf{A} \in \mathbb{R}^{n \times n}$: the RREF of $[\mathbf{A} \mid \mathbf{I}_n]$ will be $[\mathbf{I}_n \mid \mathbf{A}^{-1}]$.

2.4 Vector spaces

A real-valued **vector space** $V = (\mathcal{V}, +, \cdot)$ consists of a set \mathcal{V} and two operations

$$+ : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V} \quad \text{vector addition}$$

$$\cdot : \mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V} \quad \text{scalar multiplication}$$

where $(\mathcal{V}, +)$ is an Abelian group* with neutral element $\mathbf{0}$

$$\text{and } \forall \lambda, \psi \in \mathbb{R}, \mathbf{x}, \mathbf{y} \in \mathcal{V} : \lambda \cdot (\mathbf{x} + \mathbf{y}) = \lambda \cdot \mathbf{x} + \lambda \cdot \mathbf{y}$$

$$(\lambda + \psi) \cdot \mathbf{x} = \lambda \cdot \mathbf{x} + \psi \cdot \mathbf{x}$$

$$\lambda \cdot (\psi \cdot \mathbf{x}) = (\lambda\psi) \cdot \mathbf{x}$$

and the neutral element w.r.t. scalar multiplication is 1, such that $\forall \mathbf{x} \in \mathcal{V} : 1 \cdot \mathbf{x} = \mathbf{x}$.

* closed, associative, commutative, $\forall \mathbf{x} \in \mathcal{V} : \mathbf{x} + \mathbf{0} = \mathbf{x}$, $\forall \mathbf{x} \in \mathcal{V} \exists \mathbf{y} \in \mathcal{V} : \mathbf{x} + \mathbf{y} = \mathbf{0}$

We will denote a vector space $(\mathcal{V}, +, \cdot)$ by V , and assume $+$ and \cdot are the standard vector addition and scalar multiplication.

We'll often write $\mathbf{x} \in V$ to simplify notation.

We also often omit the dot in scalar multiplication: $\lambda \mathbf{x} = \lambda \cdot \mathbf{x}$

Vector subspaces

Let $V = (\mathcal{V}, +, \cdot)$ be a vector space, and $\mathcal{U} \subseteq \mathcal{V}$ with $\mathcal{U} \neq \emptyset$.

Then $U = (\mathcal{U}, +, \cdot)$ is a **vector subspace** of V if:

- \mathcal{U} contains the neutral element w.r.t. vector addition: $\mathbf{0} \in \mathcal{U}$
- U is closed w.r.t. vector addition: $\forall \mathbf{x}, \mathbf{y} \in \mathcal{U} : \mathbf{x} + \mathbf{y} \in \mathcal{U}$
- U is closed w.r.t. scalar multiplication: $\forall \lambda \in \mathbb{R}, \mathbf{x} \in \mathcal{U} : \lambda \mathbf{x} \in \mathcal{U}$

2.5 Linear independence

A **linear combination** of $\mathbf{x}_1, \dots, \mathbf{x}_k$ in vector space V is any vector $\mathbf{v} \in V$ of the form

$$\mathbf{v} = \lambda_1 \mathbf{x}_1 + \dots + \lambda_k \mathbf{x}_k \quad \text{with } \lambda_1, \dots, \lambda_k \in \mathbb{R}.$$

A set $\mathbf{x}_1, \dots, \mathbf{x}_k \in V$ is **linearly dependent** if there is a non-trivial linear combination

$$\lambda_1 \mathbf{x}_1 + \dots + \lambda_k \mathbf{x}_k = \mathbf{0} \quad \text{with at least one } \lambda_i \neq 0.$$

If the only way to form $\mathbf{0}$ is with $\lambda_1, \dots, \lambda_k = 0$, the set is **linearly independent**.

Linear independence implies no vector in the set can be written as a linear combination of the others (no redundancy).

In row-echelon form, non-pivot columns can be expressed as linear combinations of pivot columns on their left. So columns of \mathbf{A} are linearly independent if and only if all columns in the REF of \mathbf{A} are pivot columns.

2.6 Basis and rank

The **span** of $\mathcal{A} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ in vector space V is all possible linear combinations of the vectors in \mathcal{A} . If $\text{span}[\mathcal{A}] = V$, then \mathcal{A} is a **generating set** of V .

If the vectors in generating set \mathcal{A} are linearly independent, then \mathcal{A} is a **basis** of V .

The canonical / standard basis of \mathbb{R}^n consists of the columns of I_n .

Every basis of vector space V has the same number of vectors; the **dimension** of V .

If $U \subseteq V$, $\dim(U) \leq \dim(V)$, and $\dim(U) = \dim(V)$ if and only if $U = V$.

Finding a basis of $U = \text{span}[\mathbf{x}_1, \dots, \mathbf{x}_m]$:

1. write the spanning vectors as columns of matrix \mathbf{A}
2. determine the row-echelon form of \mathbf{A}
3. the spanning vectors associated with pivot columns are a basis of U

The **rank** of matrix \mathbf{A} , written as $\text{rk}(\mathbf{A})$, is the number of linearly independent columns (or rows) of \mathbf{A} . Note: $\text{rk}(\mathbf{A}) = \text{rk}(\mathbf{A}^T)$.

If U is the subspace spanned by the columns of \mathbf{A} , then $\dim(U) = \text{rk}(\mathbf{A})$.

Later we'll call this U the *image* or *range* of \mathbf{A} .

If W is the subspace spanned by the rows of \mathbf{A} , then $\dim(W) = \text{rk}(\mathbf{A})$.

Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is invertible if and only if $\text{rk}(\mathbf{A}) = n$.

For $\mathbf{A} \in \mathbb{R}^{m \times n}$, the subspace of solutions for $\mathbf{Ax} = \mathbf{b}$ has dimension $n - \text{rk}(\mathbf{A})$.

Later we'll call this subspace the *kernel* or *null space* of \mathbf{A} .

Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ has **full rank** if $\text{rk}(\mathbf{A}) = \min\{m, n\}$. Otherwise, \mathbf{A} is **rank deficient**.

2.7 Linear mappings

For vector spaces V and W , the mapping $\Phi : V \rightarrow W$ is a **linear mapping** if

$$\forall \lambda, \psi \in \mathbb{R}, \mathbf{x}, \mathbf{y} \in \mathcal{V} : \Phi(\lambda \mathbf{x} + \psi \mathbf{y}) = \lambda \Phi(\mathbf{x}) + \psi \Phi(\mathbf{y})$$

If Φ is bijective*, there exists an **inverse mapping** $\Psi : W \rightarrow V$ such that $\Psi(\Phi(\mathbf{x})) = \mathbf{x}$.

Identity mapping in V : $\text{id}_V : V \rightarrow V$, with $\text{id}_V(\mathbf{x}) = \mathbf{x}$.

Let $B = (\mathbf{b}_1, \dots, \mathbf{b}_n)$ be an ordered basis of vector space V .

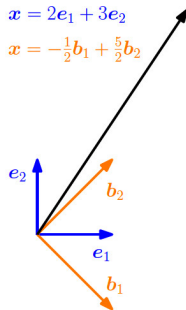
Any $\mathbf{x} \in V$ can be written as $\mathbf{x} = \alpha_1 \mathbf{b}_1 + \dots + \alpha_n \mathbf{b}_n$ and we

call $\alpha_1, \dots, \alpha_n$ the **coordinates** of \mathbf{x} w.r.t. B .

* injective: if $\Phi(\mathbf{x}) = \Phi(\mathbf{y})$ then $\mathbf{x} = \mathbf{y}$

surjective: $\Phi(V) = W$

bijjective: injective and surjective



Let $B = (\mathbf{b}_1, \dots, \mathbf{b}_n)$ and $C = (\mathbf{c}_1, \dots, \mathbf{c}_m)$ be bases of vector spaces V and W . Consider a linear mapping $\Phi : V \rightarrow W$, such that $\Phi(\mathbf{b}_j) = \alpha_{1,j}\mathbf{c}_1 + \dots + \alpha_{m,j}\mathbf{c}_m$. The matrix \mathbf{A} with elements $\alpha_{i,j}$ is the **transformation matrix** of Φ (w.r.t. B and C).

If $\hat{\mathbf{x}}$ is the coordinate vector of \mathbf{x} , and $\hat{\mathbf{y}}$ that of $\mathbf{y} = \Phi(\mathbf{x})$, then $\hat{\mathbf{y}} = \mathbf{A}\hat{\mathbf{x}}$.

Basis change

Consider two ordered bases $B = (\mathbf{b}_1, \dots, \mathbf{b}_n)$ and $\tilde{B} = (\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_n)$ of V ,
and two ordered bases $C = (\mathbf{c}_1, \dots, \mathbf{c}_m)$ and $\tilde{C} = (\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_m)$ of W ,

Let $\mathbf{A}_\Phi \in \mathbb{R}^{m \times n}$ be the transformation matrix of $\Phi : V \rightarrow W$ w.r.t. bases B and C ,
and $\tilde{\mathbf{A}}_\Phi \in \mathbb{R}^{m \times n}$ the corresponding transformation matrix w.r.t. bases \tilde{B} and \tilde{C}

Then $\tilde{\mathbf{A}}_\Phi = \mathbf{T}^{-1}\mathbf{A}_\Phi\mathbf{S}$ with $\mathbf{S} \in \mathbb{R}^{n \times n}$ the t.m. of id_V that maps coords w.r.t. \tilde{B} to B ,
and $\mathbf{T} \in \mathbb{R}^{m \times m}$ the t.m. of id_W that maps coords from \tilde{C} to C .

Image and kernel

The **image/range** of $\Phi : V \rightarrow W$ is $\text{Im}(\Phi) = \Phi(V) = \{\mathbf{w} \in W \mid \exists \mathbf{v} \in V : \Phi(\mathbf{v}) = \mathbf{w}\}$.

The **kernel/null space** of $\Phi : V \rightarrow W$ is $\text{ker}(\Phi) = \Phi^{-1}(\mathbf{0}) = \{\mathbf{v} \in V : \Phi(\mathbf{v}) = \mathbf{0}\}$.

$\text{Im}(\Phi)$ is a subspace of W , and $\text{ker}(\Phi)$ is a subspace of V .

For the mapping $\Phi(\mathbf{x}) = \mathbf{Ax}$, $\text{Im}(\Phi)$ is the column space of \mathbf{A} ($\text{span}[\text{columns of } \mathbf{A}]$), and $\text{ker}(\Phi)$ is all solutions to $\mathbf{Ax} = \mathbf{0}$.

Rank-nullity theorem: $\dim(V) = \dim(\text{Im}(\Phi)) + \dim(\text{ker}(\Phi))$

2.8 Affine subspaces

Let V be a vector space, $\mathbf{x}_0 \in V$ and $U \subseteq V$ a subspace of V . The subset L , with $L = \{\mathbf{x}_0 + \mathbf{u} : \mathbf{u} \in U\}$, is called an **affine subspace** (or linear manifold, or hyperplane).

U is the direction space, and \mathbf{x}_0 is the support point.

If $\mathbf{x}_0 \notin U$, the affine subspace is not a vector subspace because it won't contain $\mathbf{0}$.

If $(\mathbf{b}_1, \dots, \mathbf{b}_k)$ is a basis of U , then any element $\mathbf{x} \in L$ can be written as

$$\mathbf{x} = \mathbf{x}_0 + \lambda_1 \mathbf{b}_1 + \dots + \lambda_k \mathbf{b}_k.$$

An **affine mapping** from V to W has the form $\phi(\mathbf{x}) = \mathbf{a} + \Phi(\mathbf{x})$, where $\Phi : V \rightarrow W$ is a linear mapping, and \mathbf{a} is a translation vector.