

Birla Institute of Technology and Science, Pilani
Work Integrated Learning Programmes Division

Important Definitions and Concepts in * ZC416

Mathematical Foundations for Data Science

&

Mathematical Foundations for Machine Learning

G. Venkiteswaran



Contents

1 Matrix Algebra	3
2 Solutions of Linear Systems	4
3 Vector Spaces and Linear Transformations	5
4 Inner Product and Orthogonality	7
5 Eigenvalues and Eigenvectors	10
6 Matrix Decompositions	11
7 Calculus and Vector Calculus	12
8 Gradient Descent Methods	15
9 Principal Component Analysis	16
10 KKT Conditions and Strong Duality	18
11 Support Vector Machine	19

Important Definitions and Properties

1 Matrix Algebra

1. A **matrix** is a rectangular array of numbers or functions.
2. The **size of a matrix** is defined as $m \times n$ where m is the number of rows and n is the number of columns. A matrix with $m = n$ is called a square matrix. The element in the i^{th} row and j^{th} column is denoted by a_{ij} .
3. A **zero matrix** is a matrix in which all elements are zero and denoted by $\mathbf{0}$.
4. **Equality of two matrices** is established if and only if the sizes are the same and the corresponding entries are equal.
5. **Addition of two matrices** $\mathbf{A}(= a_{ij})$ and $\mathbf{B}(= b_{ij})$ of the same size is the matrix with elements $a_{ij} + b_{ij}$.
6. **Scalar multiplication** of a matrix $\mathbf{A}(= a_{ij})$ with a scalar c is a matrix whose elements are ca_{ij} .
7. **Matrix multiplication** of $\mathbf{A}_{m \times n}(= a_{ij})$ and $\mathbf{B}_{n \times p}(= b_{ij})$ yields a matrix $\mathbf{C}_{m \times p}$ whose ij^{th} element is $\sum_{k=1}^n a_{ik}b_{kj}$. The equality of the number of columns of \mathbf{A} and the number of rows of \mathbf{B} should be noted. Also, in general $\mathbf{AB} \neq \mathbf{BA}$.
8. **Properties of matrices** - assuming suitable sizes for matrices \mathbf{A}, \mathbf{B} and \mathbf{C} and c, k being scalars.
 - i) **Associativity under addition:** $(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$
 - ii) **Commutativity under addition:** $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$
 - iii) **Distributivity:** $c(\mathbf{A} + \mathbf{B}) = c\mathbf{A} + c\mathbf{B}$ and $c(k\mathbf{A}) = ck(\mathbf{A})$
 - iv) **Additive identity:** \exists a matrix $\mathbf{0}$ such that $\mathbf{A} + \mathbf{0} = \mathbf{A}$
 - v) **Additive inverse:** \exists a matrix $-\mathbf{A}$ such that $\mathbf{A} + (-\mathbf{A}) = \mathbf{0}$
9. **Transpose** of a matrix $\mathbf{A}(= a_{ij})$ is obtained by changing the rows to columns and columns to rows and denoted by $\mathbf{A}^T(= a_{ji}^T)$. By definition $a_{ij}^T = a_{ji}$. The important properties of transposes are summarized below.
 - i) $(\mathbf{A}^T)^T = \mathbf{A}$
 - ii) $(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$

- iii) $(c\mathbf{A})^T = c\mathbf{A}^T$, for any scalar c
- iv) $(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T$

10. Types of matrices

- i) **Symmetric** if $\mathbf{A}^T = \mathbf{A}$
- ii) **Skew-symmetric** if $\mathbf{A}^T = -\mathbf{A}$
- iii) **Upper triangular** if $a_{ij} = 0$ for $i < j$
- iv) **Lower triangular** if $a_{ij} = 0$ for $i > j$
- v) **Diagonal** if $a_{ij} = 0$ for $i \neq j$
- vi) **Identity** if $a_{ij} = 1$ for $i = j$ and 0 otherwise
- vii) **Positive definite** if $\forall \mathbf{x} \neq \mathbf{0}, \mathbf{x}^T \mathbf{A} \mathbf{x} > 0$
- viii) **Positive semi-definite** if $\forall \mathbf{x} \neq \mathbf{0}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$

2 Solutions of Linear Systems

1. Elementary row operations consist of

- i) Interchanging two rows
- ii) Multiplying one row with a non-zero constant
- iii) Adding a constant multiple of one row to another

2. A matrix is said to be in **Row Echelon Form (REF)** if the following conditions are satisfied after performing the necessary elementary row operations.

- i) All elements below the leading non-zero entry in a row are zero
- ii) The leading non-zero entry in a row occurs to the right of the leading non-zero entry in the row above it
- iii) The zero rows occur below the non-zero rows.

It is said to be in **Reduced Row Echelon Form (RREF)** if apart from the above, the leading non-zero entries are scaled to 1 and all elements above the leading non-zero entries are zero. The leading non-zero entries are called the **pivots** and the columns containing them are called the **pivot columns**.

- 3. The number of non-zero rows in the REF / RREF is called the **rank** of the matrix.
- 4. A **system of linear equations** is written as $\mathbf{Ax} = \mathbf{b}$ where \mathbf{A} is a matrix of size $m \times n$. It is said to be **non-homogeneous** when $\mathbf{b} \neq \mathbf{0}$ and **homogeneous** otherwise.

5. REF applied to the augmented matrix $(A|b)$ associated with a system of linear equations $Ax = b$, where A is a square matrix is called the **forward elimination process of Gaussian elimination**. Solving the resulting system is called the **backward substitution**.
6. For the general case where A is of size $m \times n$ having rank r , we have the following scenarios
 - i) if $\text{rank}(A|b) = r$, then the system has atleast one solution. Else no solution.
 - ii) if $n > r$, then $n - r$ variables can be given arbitrary values and the values of r variables can be uniquely obtained.
7. The particular solution can be uniquely found using the pivot columns and the right hand side and the non-pivotal columns can be used to find the solution of $Ax = 0$. The general solution is then the combination of the particular solution and the general solution of $Ax = 0$.

3 Vector Spaces and Linear Transformations

1. A **binary operator** $*$ on a non-empty set X gives us a rule to perform an operation on two given elements of the set. A binary operator is said to follow the closure property if $\forall a, b \in X, a * b \in X$.
2. Let G be a non-empty set with a binary operator $*$. $\langle G, * \rangle$ is said to be a **Group** if
 - G1. $*$ is closed
 - G2. $*$ is associative, that is, $a * (b * c) = (a * b) * c \forall a, b, c \in G$
 - G3. $*$ has an identity, that is, $\forall a \in G, \exists e \in G$ such that $a * e = a$
 - G4. $*$ has an inverse, that is, $\forall a \in G, \exists b \in G$ such that $a * b = e$
3. A group $\langle G, * \rangle$ is said to be **Abelian** if $a * b = b * a \forall a, b \in G$.
4. A **Field** is an Abelian group with respect to the usual addition and multiplication.
5. Let V be a non-empty set over a field F . V is called a vector space if the following conditions are satisfied.
 - A1. Associativity for addition: $\forall u, v, \omega \in V, u + (v + \omega) = (u + v) + \omega$
 - A2. Commutativity for addition: $\forall u, v \in V, u + v = v + u$
 - A3. Additive identity: $\exists \mathbf{0} \in V$ such that $\forall v \in V, v + \mathbf{0} = v$
 - A4. Additive inverse: $\forall v \in V, \exists -v \in V$ such that $v + (-v) = \mathbf{0}$

M1. Left distributivity: $\forall \mathbf{u}, \mathbf{v} \in \mathbf{V}$ and $c \in \mathbf{F}$ $c(\mathbf{u} + \mathbf{v}) = c\mathbf{u} + c\mathbf{v}$

M2. Right distributivity: $\forall \mathbf{v} \in \mathbf{V}$ and $c, d \in \mathbf{F}$ $(c+d)\mathbf{v} = c\mathbf{v} + d\mathbf{v}$

M3. Scalar multiplication: $\forall \mathbf{v} \in \mathbf{V}$ and $c, d \in \mathbf{F}$ $c(d\mathbf{v}) = (cd)\mathbf{v}$

M4. Multiplicative identity: $\exists 1 \in \mathbf{F}$ such that $\forall \mathbf{v} \in \mathbf{V}$, $1 \cdot \mathbf{v} = \mathbf{v}$

6. The **linear span** of a set of vectors $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ is the set

$$LS(S) = \left\{ \sum_{i=1}^n \alpha_i \mathbf{v}_i \text{ where } \alpha_i \in \mathbf{F} \forall i = 1, 2, \dots, n \right\}.$$

7. A set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ is **linearly independent** if $\sum_{i=1}^n \alpha_i \mathbf{v}_i = \mathbf{0}$ has $\alpha_i = 0$ as the only solution $\forall i = 1, 2, \dots, n$. If at least one of the $\alpha_i \neq 0$, then it is said to be **linearly dependent**.

8. A set of vectors $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ is called a **basis** of a vector space V over F if S is linearly independent and spans V . An equivalent definition is that S is the maximum number of linearly independent elements in V . A vector space can have multiple bases.

9. $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^2$ are linearly dependent if there exists a $\lambda \in \mathbb{R}$ such that $\mathbf{v}_2 = \lambda \mathbf{v}_1$ or $\mathbf{v}_1 = \lambda \mathbf{v}_2$. In \mathbb{R}^2 , if we cannot express one of the two vectors as a linear multiple of the other then they are linearly independent.

10. The number of elements in a basis is called the **dimension** of V

11. Let \mathbf{V} be a vector space over F and $\mathbf{W} \subset \mathbf{V}$. \mathbf{W} is a **subspace** over F if a) $\mathbf{0} \in \mathbf{W}$ and b) $\alpha \omega_1 + \beta \omega_2 \in \mathbf{W} \forall \omega_1, \omega_2 \in \mathbf{W}$ and $\alpha, \beta \in F$.

12. **Row space** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the linear span of the row vectors of the matrix \mathbf{A} , denoted by $\text{row}(\mathbf{A})$ and is a subspace of \mathbb{R}^n .

13. **Column space** of $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the linear span of column vectors of the matrix \mathbf{A} , denoted by $\text{col}(\mathbf{A})$ and is a subspace of \mathbb{R}^m .

14. **Equivalence of dimensions** $\dim \text{row}(\mathbf{A}) = \dim \text{col}(\mathbf{A}) = \text{rank}(\mathbf{A})$.

15. The **null space** of $\mathbf{A} \in \mathbb{R}^{m \times n}$ is defined as $N(\mathbf{A}) = \{x \in \mathbb{R}^n | \mathbf{A}x = \mathbf{0}\}$ and is a subspace of \mathbb{R}^n .

16. **Rank Nullity Theorem:** for $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A}) + \dim N(\mathbf{A}) = n$.

17. A mapping $T : \mathbf{V} \rightarrow \mathbf{W}$, where \mathbf{V} and \mathbf{W} are vector spaces over the same field F is called a **linear transformation** if

- i) $T(\mathbf{v}_1 + \mathbf{v}_2) = T(\mathbf{v}_1) + T(\mathbf{v}_2) \forall \mathbf{v}_1, \mathbf{v}_2 \in \mathbf{V}$

- ii) $T(c\mathbf{v}) = cT(\mathbf{v}) \forall c \in F, \mathbf{v} \in \mathbf{V}$
- 18. The sets $\mathbf{R}(T) = \{\mathbf{u} \in \mathbf{W} \mid \mathbf{u} = T(\mathbf{v}) \text{ for some } \mathbf{v} \in \mathbf{V}\}$ and $\mathbf{N}(T) = \{\mathbf{v} \in \mathbf{V} \mid T(\mathbf{v}) = \mathbf{0}\}$ are subspaces of \mathbf{W} and \mathbf{V} and are called the **range** and **null space** respectively.

- 19. The **Rank-Nullity theorem** for linear transformation states that

$$\dim(R(T)) + \dim(N(T)) = \dim(\mathbf{V}).$$

- 20. Associated with every linear transformation $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is an $n \times m$ matrix which is the **matrix representation of the linear transformation**.
- 21. A linear transformation of a vector in \mathbb{R}^n is a combination of rotation and scaling.

4 Inner Product and Orthogonality

- 1. For \mathbf{a} and \mathbf{b} in \mathbb{R}^n , $\mathbf{a}^T \mathbf{b}$ is called the **dot product** of \mathbf{a} and \mathbf{b} and is denoted by $\langle \mathbf{a}, \mathbf{b} \rangle$ or $\mathbf{a} \cdot \mathbf{b}$.

- 2. Properties of dot product

- i) $\langle k\mathbf{u} + l\mathbf{v}, \mathbf{w} \rangle = k\langle \mathbf{u}, \mathbf{w} \rangle + l\langle \mathbf{v}, \mathbf{w} \rangle, \forall k, l \in \mathbb{R}, \forall \mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ (**linearity**)
- ii) $\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{u} \rangle \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ (**symmetry**)
- iii) $\langle \mathbf{u}, \mathbf{u} \rangle \geq 0, \forall \mathbf{u} \in \mathbb{R}^n$ $\langle \mathbf{u}, \mathbf{u} \rangle = 0$ if and only if $\mathbf{u} = \mathbf{0}$ (**positive definite**)

- 3. The **norm** of a vector $\mathbf{a} \in \mathbb{R}^n$ is defined as $\|\mathbf{a}\| = \sqrt{\langle \mathbf{a}, \mathbf{a} \rangle} \geq 0$.

- 4. Properties of norm

- i) $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \|\mathbf{a}\| \|\mathbf{b}\|$ (**Cauchy Schwarz inequality**)
- ii) $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$ (**Triangle inequality**)

- 5. For any vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$

$$-1 \leq \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\| \|\mathbf{b}\|} \leq 1$$

- 6. A bilinear mapping Ω is a mapping with two arguments and is linear in both arguments: Let \mathbf{V} be a vector space such that $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{V}$, and let $\lambda, \psi \in \mathbb{R}$. Then we have $\Omega(\lambda\mathbf{x} + \psi\mathbf{y}, \mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{z}) + \psi\Omega(\mathbf{y}, \mathbf{z})$, and $\Omega(\mathbf{x}, \lambda\mathbf{y} + \psi\mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{y}) + \psi\Omega(\mathbf{x}, \mathbf{z})$.

7. Let \mathbf{V} be a vector space and $\Omega : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ be a bilinear mapping that takes two vectors as arguments and returns a real number. Then Ω is called symmetric if $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$. Also Ω is called positive-definite if $\forall \mathbf{x} \in \mathbf{V} \setminus \{\mathbf{0}\}$, $\Omega(\mathbf{x}, \mathbf{x}) > 0$ and $\Omega(\mathbf{0}, \mathbf{0}) = 0$.
8. A bilinear mapping Ω is a mapping with two arguments and is linear in both arguments: Let \mathbf{V} be a vector space such that $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$, and let $\lambda, \psi \in \mathbb{R}$. Then we have $\Omega(\lambda\mathbf{x} + \psi\mathbf{y}, \mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{z}) + \psi\Omega(\mathbf{y}, \mathbf{z})$, and $\Omega(\mathbf{x}, \lambda\mathbf{y} + \psi\mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{y}) + \psi\Omega(\mathbf{x}, \mathbf{z})$.
9. Let \mathbf{V} be a vector space and $\Omega : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ be a bilinear mapping that takes two vectors as arguments and returns a real number. Then Ω is called symmetric if $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$. Also Ω is called positive-definite if $\forall \mathbf{x} \in \mathbf{V} \setminus \{\mathbf{0}\}$, $\Omega(\mathbf{x}, \mathbf{x}) > 0$ and $\Omega(\mathbf{0}, \mathbf{0}) = 0$.
10. A positive-definite, symmetric bilinear mapping $\Omega : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ is called an inner product. To denote an inner product on \mathbf{V} we generally write $\langle \mathbf{x}, \mathbf{y} \rangle$. The pair $(\mathbf{V}, \langle \cdot, \cdot \rangle)$ is called an inner product space.
11. For a real-valued, finite-dimensional vector space \mathbf{V} and an ordered basis B of \mathbf{V} , it holds that $\langle \cdot, \cdot \rangle : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ is an inner product if and only if there exists a symmetric, positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\langle \mathbf{x}, \mathbf{y} \rangle = \hat{\mathbf{x}}^T \mathbf{A} \hat{\mathbf{y}}$.
12. Inner products and norms are closely related in the sense that any inner product induces a norm: $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
13. Not every norm is induced by an inner product, for example the Manhattan norm.
14. For an inner product vector space $(V, \langle \cdot, \cdot \rangle)$, the induced norm $\|\cdot\|$ satisfies the Cauchy-Schwarz inequality: $\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \|\mathbf{y}\|$.
15. The **angle** between two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ is
$$\alpha = \cos^{-1} \left(\frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\| \|\mathbf{b}\|} \right).$$
16. Two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are said to be **orthogonal** if $\langle \mathbf{a}, \mathbf{b} \rangle = 0$.
17. A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is called as an **orthogonal set** if \mathbf{v}_i is orthogonal to \mathbf{v}_j , $\forall i \neq j$.
18. A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is called as an **orthonormal set** if \mathbf{v}_i is orthogonal to \mathbf{v}_j , $\forall i \neq j$ and each \mathbf{v}_i is of unit norm.
19. The **projection** of \mathbf{v}_2 onto the vector \mathbf{v}_1 is $\mathbf{v} = \lambda \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}$ where $\lambda = \frac{\langle \mathbf{v}_2, \mathbf{v}_1 \rangle}{\|\mathbf{v}_1\|}$.

20. The canonical basis or the standard basis forms an orthonormal set.
21. A matrix \mathbf{A} is said to be an orthonormal matrix if the column vectors of the matrix form a orthonormal set and $\mathbf{A}\mathbf{A}^T = \mathbf{I} = \mathbf{A}^T\mathbf{A}$.
22. **Gram-Schmidt Process** If $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ be a basis of a subspace U of \mathbb{R}^n , and

$$\begin{aligned}\mathbf{v}_1 &= \mathbf{u}_1 \\ \mathbf{v}_2 &= \mathbf{u}_2 - \frac{\mathbf{u}_2 \cdot \mathbf{v}_1}{\mathbf{v}_1 \cdot \mathbf{v}_1} \mathbf{v}_1 \\ \mathbf{v}_3 &= \mathbf{u}_3 - \frac{\mathbf{u}_3 \cdot \mathbf{v}_1}{\mathbf{v}_1 \cdot \mathbf{v}_1} \mathbf{v}_1 - \frac{\mathbf{u}_3 \cdot \mathbf{v}_2}{\mathbf{v}_2 \cdot \mathbf{v}_2} \mathbf{v}_2 \\ &\vdots \\ \mathbf{v}_m &= \mathbf{u}_m - \frac{\mathbf{u}_m \cdot \mathbf{v}_1}{\mathbf{v}_1 \cdot \mathbf{v}_1} \mathbf{v}_1 - \frac{\mathbf{u}_m \cdot \mathbf{v}_2}{\mathbf{v}_2 \cdot \mathbf{v}_2} \mathbf{v}_2 - \dots - \frac{\mathbf{u}_m \cdot \mathbf{v}_{m-1}}{\mathbf{v}_{m-1} \cdot \mathbf{v}_{m-1}} \mathbf{v}_{m-1}\end{aligned}$$

Then, $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ is an orthogonal basis for U and
 $\text{LS}\{\mathbf{v}_1, \dots, \mathbf{v}_k\} = \text{LS}\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ for $1 \leq k \leq m$.

23. Consider an inner product space $(V, \langle \cdot, \cdot \rangle)$. Define $d(\mathbf{x}, \mathbf{y})$ the distance between two vectors \mathbf{x} and \mathbf{y} to be $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle}$.
24. If we use the dot product as the inner product, then the distance is called the Euclidean distance.
25. The mapping $d : V \times V \rightarrow \mathbb{R}$ is called a metric.
26. d is positive-definite which means $d(\mathbf{x}, \mathbf{y}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in V$. $d(\mathbf{x}, \mathbf{y}) = 0 \implies \mathbf{x} = \mathbf{y}$.
27. d is symmetric which means $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in V$.
28. d obeys the triangle inequality as follows: $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}) \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in V$
29. Performing Gaussian elimination on the matrix $\mathbf{A}^T\mathbf{A}$ where \mathbf{A} contains the basis vectors as its columns. Upon Gaussian elimination on the augmented matrix we reduce $[\mathbf{A}^T\mathbf{A} | \mathbf{A}^T]$ to get $[\mathbf{U} | \mathbf{L}^{-1}\mathbf{A}^T]$ where $\mathbf{A}^T\mathbf{A} = \mathbf{L}\mathbf{U}$. $\mathbf{Q}^T = \mathbf{L}^{-1}\mathbf{A}^T$ is an orthogonal matrix whose rows are orthogonal. This is the alternautive form of **Gram-Schmidt Orthogonalization**.

5 Eigenvalues and Eigenvectors

1. The **minor** of an element a_{ij} , denoted by m_{ij} , of a square matrix $\mathbf{A} (= a_{ij})$ of size $n \times n$ is the determinant of the submatrix of \mathbf{A} obtained by deleting the i^{th} row and the j^{th} column from \mathbf{A} .
2. The **cofactor** of the element a_{ij} is defined as $(-1)^{i+j}$.
3. The **determinant** of a matrix is then defined as $\text{Det}(\mathbf{A}) = \sum_{k=1}^n a_{rk} c_{rk}$ where we have taken the product over the r^{th} row. The value is invariant for the product taken over columns. The following are some of the interesting properties.
 - i) $\text{Det}(\mathbf{A}) = 0$ whenever there is linear dependence in the rows / columns.
 - ii) Multiplying the r^{th} row by a scalar c is same as multiplying the determinant value with c .
 - iii) $\text{Det}(c\mathbf{A}) = c^n \text{Det}(\mathbf{A})$
 - iv) $\text{Det}(\mathbf{AB}) = \text{Det}(\mathbf{A})\text{det}(\mathbf{B})$
4. The **adjoint** of a matrix \mathbf{A} , denoted by $\text{adj}(\mathbf{A})$ is defined as the transpose of the cofactor matrix.
5. The **inverse** of a matrix \mathbf{A} is defined as $\mathbf{A}^{-1} = \frac{\text{adj}(\mathbf{A})}{\text{Det}(\mathbf{A})}$.
6. An alternative way to calculate the inverse of \mathbf{A} is to start with the augmented matrix $[\mathbf{A} | \mathbf{I}]$ and use elementary row operations to convert this to $[\mathbf{I} | \mathbf{A}^{-1}]$. This procedure is called the **Gauss Jordan** method.
7. The roots of the **characteristic equation** $\text{Det}(\mathbf{A} - \lambda \mathbf{I}) = 0$ are called the **eigenvalues**. Any non-zero vector \mathbf{x} which satisfies the equation $\mathbf{Ax} = \lambda \mathbf{x}$ is called the **eigenvector** corresponding to the eigenvalue λ .
8. There are n eigenvalues for a matrix \mathbf{A} of size $n \times n$ and they can be real or complex.
9. For a square matrix the **rank** is the number of non-zero eigenvalues.
10. A symmetric positive definite matrix has full rank.
11. The eigenvalues of a positive definite matrix are always positive.

12. The sum of diagonal elements of a matrix, called the **trace** is equal to the sum of the eigenvalues and the product of the eigenvalues is equal to the determinant.
13. The eigenvalues of a symmetric matrix are all real whereas that of a skew-symmetric matrix are either purely complex or zero.
14. The eigenvectors corresponding to distinct eigenvalues are linearly independent.
15. The **Spectral theorem** states that for a symmetric matrix \mathbf{A} of size $n \times n$, the eigenvalues are real and that the eigenvectors for an orthogonal basis of \mathbb{R}^n .

6 Matrix Decompositions

1. If x_1, x_2, \dots, x_n are the linearly independent eigenvectors of a matrix corresponding to eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, then

$$\mathbf{P}^{-1}\mathbf{AP} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

where \mathbf{P} is the matrix with eigenvectors as columns. This is called the **eigenvalue decomposition**.

2. Under certain conditions a matrix \mathbf{A} can be decomposed into a product of a lower triangular matrix (\mathbf{L}) and an upper triangular matrix (\mathbf{U}), that is $\mathbf{A} = \mathbf{LU}$. This method is computationally effective in solving systems of equations with the same matrix \mathbf{A} and different right hand sides. This procedure is called **LU decomposition method**.
3. The steps in LU decomposition include, solving for \mathbf{y} in $\mathbf{Ly} = \mathbf{b}$ and then $\mathbf{Ux} = \mathbf{y}$. The following methods are popular.
 - i) **Doolittle's method** if $L_{ii} = 1 \forall i$
 - ii) **Crout's method** if $U_{ii} = 1 \forall i$
 - iii) **Cholesky's method** if \mathbf{A} is positive definite and $\mathbf{L} = \mathbf{U}^T$.
4. Any matrix $\mathbf{A}_{m \times n}$ can be written as $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ where $\mathbf{U}_{m \times m}, \mathbf{V}_{n \times n}$ are orthonormal matrices and $\Sigma_{m \times n}$ is the matrix of singular values. This is called the **Singular Value Decomposition** of \mathbf{A} .
5. The matrix \mathbf{U} consists of eigenvectors of \mathbf{AA}^T written as columns and \mathbf{V} has the eigenvectors of $\mathbf{A}^T\mathbf{A}$ written as columns.
6. The elements of Σ are all non-negative and are arranged in decreasing order along the diagonal.

7. If $\mathbf{u}_1, \mathbf{v}_2, \dots, \mathbf{u}_m$ and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are the columns of \mathbf{U} and \mathbf{V} and σ_{ii} , $i = 1, 2, \dots, r$ are the non-zero diagonal elements of Σ , of a matrix \mathbf{A} of rank r , then a rank k approximation of \mathbf{A} (where $k \leq r$) is given by $\sum_{i=1}^k \mathbf{u}_i \sigma_{ii} \mathbf{v}_i^T$.
8. The first diagonal element of Σ is called the 2 norm of \mathbf{A} .
9. If \mathbf{A} is an $m \times n$ matrix with $m \geq n$ and all the columns linearly independent, then \mathbf{A} has a decomposition of the form $\mathbf{Q}\mathbf{R}$ where \mathbf{Q} is a matrix whose columns are orthogonal and \mathbf{R} is an upper triangular matrix. This is called the **QR decomposition** of \mathbf{A} . The basic idea is to use Gram-Schmidt Orthogonalization process to get \mathbf{Q} and get \mathbf{R} via $\mathbf{Q}^T\mathbf{A}$.

7 Calculus and Vector Calculus

1. The **derivative** of f at x is defined as the limit

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

2. The **Taylor polynomial** of degree n of $f : \mathbb{R} \rightarrow \mathbb{R}$ at x_0 is defined as

$$T_n(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k$$

3. The **Taylor series** of smooth (continuously differentiable infinite many times) function $f : \mathbb{R} \rightarrow \mathbb{R}$ at x_0 is defined as

$$T_\infty(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k$$

4. Rules for differentiation. We denote the derivative of f by f'
 - i) **Product Rule:** $(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$
 - ii) **Sum Rule:** $(f(X) + g(X))' = f'(X) + g'(X)$
 - iii) **Quotient Rule:** $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$
 - iv) **Chain Rule:** $(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$

5. For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $x \rightarrow f(x)$, $x \in \mathbb{R}^n$ of n variables x_1, \dots, x_n we define the **partial derivatives** as

$$\begin{aligned}\frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h} \\ \frac{\partial f}{\partial x_2} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2 + h, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h}\end{aligned}$$

We collect them in the row vector called the gradient of f or **Jacobian**

$$\Delta_x f = \mathbf{grad} f = \frac{df}{dx} = \left[\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_n} \right]$$

6. The rules for partial differentiation are

- i) **Product rule:** $\frac{\partial}{\partial x}(f(x)g(x)) = \frac{\partial f}{\partial x}g(x) + f(x)\frac{\partial g}{\partial x}$
- ii) **Sum rule:** $\frac{\partial}{\partial x}(f(x) + g(x)) = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$
- iii) **Chain rule:** $\frac{\partial}{\partial x}(g \circ f)(x) = \frac{\partial}{\partial x}(g(f(x))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x}$

7. To compute the gradient of f with respect to t , we need to apply the chain rule for multivariate functions as

$$\frac{df}{dt} = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] \left[\begin{array}{c} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{array} \right] = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

where d denotes the gradient and ∂ partial derivatives.

8. For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a vector $x = [x_1, \dots, x_n]^T$ corresponding vector of function values is given as

$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix} \in \mathbb{R}^m$$

where each $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$

9. **Taylor's Theorem:** Suppose $f : (a, b) \rightarrow R$ is a function on (a, b) , where a, b in R with $a < b$. Assume that f is n -times differentiable in the open interval (a, b) and $f, f', f'', \dots, f^{n-1}$ all extend continuously to the closed interval $[a, b]$, such that the extended functions are still called $f, f', f'', \dots, f^{n-1}$. Then there exists $c \in (a, b)$ such that

$$f(b) = \sum_{k=0}^{k=n-1} \frac{f^k(a)}{k!} (b-a)^k + \frac{f^n(c)}{n!} (b-a)^n$$

10. For $n = 1$, the statement of Taylor's theorem boils down to the **Mean-Value Theorem** which is that if a function f is continuous on $[a, b]$ and differentiable on the interval (a, b) , then there exists a value $c \in (a, b)$ such that $f'(c) = \frac{f(b)-f(a)}{b-a}$

11. The **Hessian** is the collection of all second-order partial derivatives. If $f(x, y)$ is a twice (continuously) differentiable function, then $\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x}$ i.e., the order of differentiation does not matter, and the corresponding Hessian matrix is symmetric. The Hessian is denoted as

$$\nabla_{x,y}^2 f(x, y) = \begin{bmatrix} f_{xx} & f_{xy} \\ f_{xy} & f_{yy} \end{bmatrix}$$

12. For a function $f(x, y)$ which is twice differentiable in a neighbourhood of the point (a, b) and $f_x(a, b) = f_y(a, b) = 0$, the expression for $f(a + h, b + k)$ simplifies to $f(a + h, b + k) - f(a, b) = \frac{1}{2}(h^2 f_{xx} + 2hk f_{xy} + k^2 f_{yy})|_{a+ch,b+ck} =: Q(c)$. If $Q(0) \neq 0$, the sign of $Q(c)$ for small c will be the same as the sign of $Q(0)$ for sufficiently small values of h and k .

13. Since $f_{xx}Q(0) = (hf_{xx} + kf_{xy})^2 + (f_{xx}f_{yy} - f_{xy}^2)k^2$.

- i) If $f_{xx} < 0$ and $f_{xx}f_{yy} - f_{xy}^2 > 0$ at (a, b) then $Q(0) < 0$ for all sufficiently small non-zero values of h and k , then f has a local maximum value at (a, b) .
- ii) If $f_{xx} > 0$ and $f_{xx}f_{yy} - f_{xy}^2 > 0$ at (a, b) then $Q(0) > 0$ for all sufficiently small non-zero values of h and k , then f has a local minimum value at (a, b) .
- iii) If $f_{xx}f_{yy} - f_{xy}^2 < 0$ at (a, b) there are combinations of small values for h and k for which $Q(0) > 0$ and other combinations of h and k for which $Q(0) < 0$. This means that f has a saddle point at (a, b) .
- iv) If $f_{xx}f_{yy} - f_{xy}^2 = 0$ at (a, b) another test is needed.

14. The gradient of an $m \times n$ matrix \mathbf{A} with respect to a $p \times q$ matrix \mathbf{B} , the resulting Jacobian would be an $(m \times n) \times (p \times q)$, i.e., a four-dimensional **tensor** J , whose entries are given as

$$J_{ijkl} = \frac{\partial \mathbf{A}_{ij}}{\partial \mathbf{B}_{kl}}$$

15. Some useful **gradient identities**

i) $\frac{\partial}{\partial X} f(X)^T = \left(\frac{\partial f(X)}{\partial X} \right)^T$

- ii) $\frac{\partial}{\partial X} \text{tr}(f(X)) = \text{tr}\left(\frac{\partial f(X)}{\partial X}\right)$
- iii) $\frac{\partial}{\partial X} \det(f(X)) = \det(f(x)) \text{tr}\left(f(X)^{-1} \frac{\partial f(X)}{\partial X}\right)$
- iv) $\frac{\partial}{\partial X} f(X)^{-1} = -f(X)^{-1} \frac{\partial f(X)}{\partial X} f(X)^{-1}$
- v) $\frac{\partial a^T X^{-1} b}{\partial X} = -(X^{-1})^T a b^T (X^{-1})^T$
- vi) $\frac{\partial x^T a}{\partial x} = a^T$
- vii) $\frac{\partial a^T x}{\partial x} = a^T$
- viii) $\frac{\partial a^T X b}{\partial X} = ab^T$
- ix) $\frac{\partial x^T B}{\partial x} = x^T (B + B^T)$
- x) $\frac{\partial}{\partial s}(x - As)^T W(x - As) = -2(x - As)^T W A$, for symmetric W .

8 Gradient Descent Methods

1. Gradient descent method: find the optimum of J , say at $J(\mathbf{x}_*)$, we can start at some initial point \mathbf{x}_0 and then iterate according to $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha_i ((\nabla J)(\mathbf{x}_i))^T$, where α is the **learning rate**.
2. The standard gradient descent procedure is a batch optimization method in that the update step considers the gradient of the entire loss function $L(\theta)$, i.e $\theta_{i+1} = \theta_i - \alpha_i \nabla L(\theta_i)^T = \theta_i - \alpha_i \sum_{n=1}^{N} \nabla L_n(\theta_i)^T$.
3. With a learning rate dependent on time, the update step becomes $\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha_t \nabla J$.
4. In line search, the step-size α_t is computed as $\alpha_t = \min_{\alpha} J(\mathbf{w}_t + \alpha \mathbf{g}_t)$.
5. The first step in optimization is to identify a range $[a, b] = [0, \alpha_{\max}]$ in which to perform the search for the optimum α .
6. It is then possible to narrow the search interval by using **binary search** or **golden section search** methods.
 - i) In binary search, if the objective function is found to be increasing at $\frac{a+b}{2}$, we narrow the interval to $[a, \frac{a+b+\epsilon}{2}]$ and continue the search. Otherwise we narrow the interval to $[\frac{a+b}{2}, b]$ and continue the search. ϵ is usually taken as 10^{-8} .

- ii) In golden section search, When $\alpha = a$ yields the minimum for the objective function, i.e $H(\alpha)$, we can drop the interval $(m_1, b]$. Similarly when $\alpha = b$ yields the minimum for $H(\alpha)$ we can drop the interval $[a, m_2]$. When $\alpha = m_1$ is the value at which the minimum is achieved we can drop $(m_2, b]$. When $\alpha = m_2$ is the value at which the minimum is achieved we can drop $[a, m_1]$.
- 7. In case of **mean centering** a vector of column-wise means is subtracted from each data point.
- 8. In case of **feature normalization**, each feature value is divided by its standard deviation.
- 9. In case of **min-max normalization** we scale the j th feature of the i th datapoint as follows: $x_{ij} = \frac{x_{ij} - \min_j}{\max_j - \min_j}$.
- 10. Modified Gradient Descent Methods
 - i) The set S of data points can be treated as a sample and a sample-centric objective function can be constructed as follows: $J(S) = \sum_{i \in S} (\mathbf{w}^T \mathbf{X}_i - y_i)^2$. This is called as **mini batch gradient descent**.
 - ii) In the extreme case S can contain only one index chosen at random, the approach is called as **stochastic gradient descent**.
 - iii) The normal update procedure for gradient descent can be written as $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$ where $\mathbf{v} \leftarrow -\alpha \frac{\partial J}{\partial \mathbf{w}}$.
 - iv) The gradient descent with **momentum based**, for $\beta \in (0, 1)$, the update can be written as $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$ where $\mathbf{v} \leftarrow \beta \mathbf{v} - \alpha \frac{\partial J}{\partial \mathbf{w}}$.
 - v) **AdaGrad Method:** The update step becomes $w_i \leftarrow w_i - \frac{\alpha}{\sqrt{A_i}} \frac{\partial J}{\partial w_i}$, $\forall i$ where $A_i \leftarrow A_i + \left(\frac{\partial J}{\partial w_i} \right)^2$, $\forall i$.
 - vi) **RMS Prop Method:** The update step is $w_i \leftarrow w_i - \frac{\alpha}{\sqrt{A_i}} \frac{\partial J}{\partial w_i}$, $\forall i$ where $A_i \leftarrow \rho A_i + (1 - \rho) \left(\frac{\partial J}{\partial w_i} \right)^2$ with $\rho \in (0, 1)$.
 - vii) **Adams Method:** The following update is used at the t th iteration: $w_i \leftarrow w_i - \frac{\alpha_t F_i}{\sqrt{A_i}}$ where $\alpha_t = \alpha \frac{\sqrt{1 - \rho^t}}{1 - \rho_f^t}$. And $A_i \leftarrow \rho A_i + (1 - \rho) \left(\frac{\partial J}{\partial w_i} \right)^2$ with $\rho \in (0, 1)$, $F_i \leftarrow \rho_f F_i + (1 - \rho_f) \frac{\partial J}{\partial w_i}$ with $\rho_f \in (0, 1)$.

9 Principal Component Analysis

1. Consider an iid dataset $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}, \mathbf{x}_n \in \mathbb{R}^D$ with mean $\mathbf{0}$ which possesses the **data covariance matrix** $S = \frac{1}{N} \sum_{n=1}^{n=N} \mathbf{x}_n \mathbf{x}_n^T$.

2. We assume there exists a lower-dimensional compressed representation \mathbf{z}_n of \mathbf{x}_n such that $\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n$ where the **projection matrix** $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_m] \in \mathbb{R}^{D \times M}$.
3. The columns of \mathbf{B} are orthonormal which means $\mathbf{b}_i^T \mathbf{b}_j = 0$ when $i \neq j$ and $\mathbf{b}_i^T \mathbf{b}_i = 1$.
4. There exists a linear relationship between the original data \mathbf{x} , its low-dimensional code \mathbf{z} and the compressed data $\tilde{\mathbf{x}}$: $\mathbf{z} = \mathbf{B}^T \mathbf{x}$, and $\tilde{\mathbf{x}} = \mathbf{B}\mathbf{z}$ for a suitable matrix \mathbf{B} .
5. PCA can then be viewed as a dimensionality reduction algorithm that maximizes the variance in the low-dimensional representation of the data to retain as much information as possible.
6. For the data covariance matrix $\mathbf{S} = \frac{1}{N} \sum_{n=1}^{n=N} \mathbf{x}_n \mathbf{x}_n^T$ we assume centred data, and can make this assumption without loss of generality.
7. Finding the direction \mathbf{b}_1 that maximizes variance can be set up as a constrained optimization problem

$$\begin{aligned} & \max \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 \text{ subject to} \\ & \|\mathbf{b}_1\| = 1 \end{aligned}$$

8. Our objective function boils down to maximizing λ which means we are looking for the eigenvector of \mathbf{S} that corresponds to its largest eigenvalue.
9. The m th **principal component** can be found by subtracting from the data the contribution of the first $m - 1$ components $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{m-1}$. Essentially, we are trying to find principal components that compress the remainder of the information.
10. We then arrive at a new data matrix $\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X}$ where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ contains the data points as column vectors and $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$ is a projection matrix that projects \mathbf{X} onto the subspace spanned by $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{m-1}$. $\hat{\mathbf{S}}$ is the covariance matrix of the data matrix $\hat{\mathbf{X}}$.
11. The sets of eigenvectors for $\hat{\mathbf{S}}$ and \mathbf{S} are the same.
12. Assume the SVD of X as $X = U\Sigma V^T$. Then

$$S = \frac{1}{N} XX^T = \frac{1}{N} U\Sigma\Sigma^T U^T$$

13. The eigenvalues λ_d of S are related to the singular values of X via

$$\lambda_d = \frac{\sigma_d^2}{N} \quad (1)$$

14. Consider the best rank-M approximation of X defined as \tilde{X}_M

$$\tilde{X}_M = \operatorname{argmin}_{\operatorname{rank}(A) \leq M} \|X - A\|_2 \quad (2)$$

15. The eigenvectors of XX^T can be computed from the eigenvectors of X^TX using the equation

$$\frac{1}{N}XX^TXc_m = \lambda_m Xc_m \quad (3)$$

16. Steps for the computation of the PCA:

- i) We need to **standardize** x_* using the mean and **standard deviation** of the training data in the d th dimension

$$x_*^{(d)} = \frac{x_*^{(d)} - \mu_d}{\sigma_d}, \quad d = 1, \dots, D \quad (4)$$

where $x_*^{(d)}$ is the d th component of x_* .

- ii) We obtain the projection as

$$\tilde{x} = BB^Tx_* \quad (5)$$

- iii) The coordinates are

$$z_* = B^Tx_* \quad (6)$$

with respect to the basis of the principal subspace.

10 KKT Conditions and Strong Duality

1. The **primal problem** is $\min f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0, 1 \leq i \leq m$. Optimization is performed over the primal variables \mathbf{x} .
2. We create the **Lagrangian** of the given constrained optimization problem as follows: $\mathfrak{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$, where $\lambda_i \geq 0$ for all i .
3. The associated Lagrangian **dual** problem is $\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \min_{\mathbf{x} \in \mathbb{R}^d} \mathfrak{L}(\mathbf{x}, \boldsymbol{\lambda})$ subject to $\boldsymbol{\lambda} \geq 0$ where $\boldsymbol{\lambda}$ are dual variables.
4. **Minmax inequality:** $\mathbf{x}, \mathbf{y}: \max_{\mathbf{y}} \min_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x}} \max_{\mathbf{y}} \phi(\mathbf{x}, \mathbf{y})$.
5. A set C is called a **convex set** C if for any $x, y \in C$, $\theta x + (1-\theta)y \in C$, for $0 \leq \theta \leq 1$.
6. The function is a **convex function** if for any $\mathbf{x}, \mathbf{y} \in C$, $f(\theta \mathbf{x} + (1-\theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1-\theta)f(\mathbf{y})$

7. Another way of looking at a convex function is to use the gradient: for any two points \mathbf{x} and \mathbf{y} , we have $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla_{\mathbf{x}}f(\mathbf{x})(\mathbf{y} - \mathbf{x})$.
8. For a primal optimization problem $\min f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0$ for $i = 1, 2, \dots, m$ and $h_j(\mathbf{x}) = 0$ for $j = 1, 2, \dots, p$, we say that it obeys **Slater's condition** if the objective function f is convex, the constraints g_i are all convex and the constraint functions h_i are all linear and there exists a point $\bar{\mathbf{x}}$ in the interior of the region, i.e $g_i(\bar{\mathbf{x}}) < 0$ for all $i \in [m]$ and $h_j(\bar{\mathbf{x}}) = 0$ for all $j \in [p]$.
9. Suppose Slater's condition holds and the region has a non-empty interior. Then we have **strong duality**.
10. Given a primal optimization problem, we say that \mathbf{x}^* and $(\lambda^*, \nu^*) \in \mathbb{R}^m \times \mathbb{R}^p$ respect the **Karash-Kuhn-Tucker conditions** if:
 - i) $g_i(\mathbf{x}^*) \leq 0 \forall i \in [m]$.
 - ii) $h_i(\mathbf{x}^*) = 0 \forall i \in [p]$.
 - iii) $\lambda_i^* \geq 0 \forall i \in [m]$.
 - iv) $\lambda_i^* g_i(\mathbf{x}^*) = 0 \forall i \in [m]$.
 - v) $\nabla f(\mathbf{x}^*) + \sum_{i=1}^{i=m} \lambda_i^* \nabla g_i(\mathbf{x}^*) + \sum_{i=1}^{i=p} \nu_i^* \nabla h_i(\mathbf{x}^*) = 0$.
11. For any optimization problem, if strong duality holds then any primal optimal solution \mathbf{x}^* and dual optimal solution $(\lambda^*, \nu^*) \in \mathbb{R}^m \times \mathbb{R}^p$ respect the KKT conditions. Conversely if f and g_i are convex for all $i \in [m]$ and h_i are affine for all $i \in [p]$ then the KKT conditions are sufficient for strong duality.

11 Support Vector Machine

1. $\boldsymbol{\omega}$ is a normal vector to the hyperplane $\langle \boldsymbol{\omega}, \mathbf{x} \rangle + b = 0$
2. **Hard Margin SVM**

$$\begin{aligned} & \min_{\boldsymbol{\omega}, b} \frac{1}{2} \|\boldsymbol{\omega}\|^2 \\ & \text{subject to } y_i(\langle \boldsymbol{\omega}, \mathbf{x}_i \rangle + b) \geq 1, \quad \forall i = 1, \dots, N \end{aligned}$$
3. Using the **Lagrangian** formulation and setting the partial derivatives to zero, $\boldsymbol{\omega} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$ and $\sum_{i=1}^N \alpha_i y_i = 0$.
4. **Classification** of \mathbf{x} is based on $\text{sign}(\langle \boldsymbol{\omega}, \mathbf{x} \rangle + b)$
5. **Support vectors** are those which are on the left and right margins.

6. Using KKT conditions, $\alpha_i = 0$ for \mathbf{x}_i that are not support vectors.

7. The **dual** of the Lagrangian is

$$D(\boldsymbol{\alpha}) = \max_{\boldsymbol{\alpha} \geq \mathbf{0}} \frac{1}{2} \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

8. **Soft Margin SVM**

$$\begin{aligned} & \min_{\boldsymbol{\omega}, b} \frac{1}{2} \|\boldsymbol{\omega}\|^2 + C \sum_{n=1}^N \xi_n \\ & \text{subject to } y_n (\langle \boldsymbol{\omega}, \mathbf{x}_n \rangle + b) \geq 1 - \xi_n \\ & \xi_n \geq 0, \quad n = 1, \dots, N \end{aligned}$$

9. **Hinge's Loss Function**

$$\min_{\boldsymbol{\omega}, b} \frac{1}{2} \|\boldsymbol{\omega}\|^2 + C \sum_{n=1}^N \max\{0, 1 - y_n (\langle \boldsymbol{\omega}, \mathbf{x}_n \rangle + b)\}$$

where C is the parameter that controls the width of the margin. Larger C implies lesser width and smaller C signifies a larger width.

10. A **kernel function** is a function that corresponds to an **inner product** in some expanded feature space. For a function $\phi : \mathbf{x} \rightarrow \phi(\mathbf{x})$, the dot product can be replaced by $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$. A widely used ϕ is $\phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$.

Index

- Associativity under addition, 3
- backward substitution, 5
- eigenvector, 10
- linearly dependent, 6
- orthonormal set, 8
- Abelian, 5
- AdaGrad Method, 16
- Adams Method, 16
- Addition of two matrices, 3
- Additive identity, 3
- Additive inverse, 3
- adjoint, 10
- angle, 8
- basis, 6
- binary operator, 5
- binary search, 15
- Cauchy Schwarz inequality, 7
- Chain rule, 12, 13
- characteristic equation, 10
- Cholesky's method, 11
- Classification, 19
- cofactor, 10
- Column space, 6
- Commutativity under addition, 3
- convex function, 18
- convex set, 18
- Crout's method, 11
- data covariance matrix, 16
- derivative, 12
- determinant, 10
- Diagonal, 4
- dimension, 6
- Distributivity, 3
- Doolittle's method, 11
- dot product, 7
- dual, 18, 20
- eigenvalue decomposition, 11
- eigenvalues, 10
- Elementary row operations, 4
- Equality of two matrices, 3
- Equivalence of dimensions, 6
- feature normalization, 16
- Field, 5
- forward elimination process, 5
- Gauss Jordan, 10
- Gaussian elimination, 5
- golden section search, 15
- grad, 13
- gradient identities, 14
- Gram-Schmidt Process, 9
- Group, 5
- Hard Margin SVM, 19
- Hessian, 14
- Hinge's Loss Function, 20
- homogeneous, 4
- Identity, 4
- inner product, 20
- inverse, 10
- Jacobian, 13
- Karash-Kuhn-Tucker conditions, 19
- kernel function, 20
- Lagrangian, 18, 19
- learning rate, 15
- linear span, 6
- linear transformation, 6
- linearity, 7
- linearly independent, 6
- Lower triangular, 4
- LU decomposition method, 11
- matrix, 3
- Matrix multiplication, 3

matrix representation of the linear transformation, 7
 mean centering, 16
 Mean-Value Theorem, 14
 min-max normalization, 16
 mini batch gradient descent, 16
 Minmax inequality, 18
 minor, 10
 momentum based learning, 16
 non-homogeneous, 4
 norm, 7
 null space, 6
 orthogonal, 8
 orthogonal set, 8
 partial derivatives, 13
 pivot columns, 4
 pivots, 4
 Positive definite, 4
 positive definite, 7
 Positive semi-definite, 4
 primal problem, 18
 principal component, 17
 Principal Component Analysis, 16
 Product rule, 12, 13
 projection, 8
 projection matrix, 17
 Properties of matrices, 3
 QR decomposition, 12
 Quotient rule, 12
 range space, 7
 rank, 4, 10
 Rank Nullity Theorem, 6
 Rank-Nullity theorem, 7
 Reduced Row Echelon Form (RREF), 4
 RMS Prop Method, 16
 Row Echelon Form (REF), 4
 Row space, 6
 Scalar multiplication, 3
 Singular Value Decomposition, 11
 size of a matrix, 3
 Skew-symmetric, 4
 Slater's condition, 19
 Soft Margin SVM, 20
 Spectral theorem, 11
 standard deviation, 18
 standardize, 18
 stochastic gradient descent, 16
 strong duality, 19
 subspace, 6
 Sum rule, 12, 13
 Support vectors, 19
 Symmetric, 4
 symmetry, 7
 system of linear equations, 4
 Taylor polynomial, 12
 Taylor series, 12
 Taylor's Theorem, 13
 tensor, 14
 trace, 11
 Transpose, 3
 Triangle inequality, 7
 Upper triangular, 4
 zero matrix, 3