

Machine Learning

Lecture 2 Preliminary of Algebra, Calculus, Statistics, and Probability

Hongteng Xu



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能學院
Gaoling School of Artificial Intelligence

Outline

Review

- ▶ **Data formulation:** Vector and matrix
- ▶ **Residual:** Vector and matrix norms
- ▶ **Basic operation:** Matrix-vector multiplication
- ▶ **Sample space:** Metric-measure space and its reconstruction

Outline

Review

- ▶ **Data formulation:** Vector and matrix
- ▶ **Residual:** Vector and matrix norms
- ▶ **Basic operation:** Matrix-vector multiplication
- ▶ **Sample space:** Metric-measure space and its reconstruction

Today

- ▶ Linear space and matrix analysis
 - ▶ More matrix operations and two important decompositions
- ▶ Derivation of (multi-dimensional) functions
 - ▶ Important concepts and derivations of important functions
- ▶ Statistics and probability theory
 - ▶ Connect data matrix with multi-dimensional random variables

Matrix Multiplication

$$\mathbf{B} = \mathbf{A}\mathbf{X}, \quad \mathbf{A} \in \mathbb{R}^{L \times M}, \mathbf{X} \in \mathbb{R}^{M \times N}, \mathbf{B} \in \mathbb{R}^{L \times N}. \quad (1)$$

- Element-wise representation:

$$b_{ij} = \sum_{k=1}^M a_{ik}x_{kj}, \quad \forall i = 1, \dots, L, j = 1, \dots, N \quad (2)$$

- Block-wise representation:

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \cdots & \mathbf{B}_{1n} \\ \vdots & \ddots & \vdots \\ \mathbf{B}_{l1} & \cdots & \mathbf{B}_{ln} \end{bmatrix}, \mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1m} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{l1} & \cdots & \mathbf{A}_{lm} \end{bmatrix}, \mathbf{X} = \begin{bmatrix} \mathbf{X}_{11} & \cdots & \mathbf{X}_{1n} \\ \vdots & \ddots & \vdots \\ \mathbf{X}_{m1} & \cdots & \mathbf{X}_{mn} \end{bmatrix} \quad (3)$$

$$\mathbf{B}_{ij} = \sum_{k=1}^m \mathbf{A}_{ik}\mathbf{X}_{kj}, \quad \forall i = 1, \dots, l, j = 1, \dots, n.$$

Matrix Transposition

- ▶ You can think of it as “flipping” the rows and columns

$$\begin{bmatrix} a \\ b \end{bmatrix}^T = \begin{bmatrix} a & b \end{bmatrix} \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad (4)$$

- ▶ $(\mathbf{A}^T)^T = \mathbf{A}$
- ▶ $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$
- ▶ $(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$
- ▶ $(\mathbf{A}^{-1})^T = \mathbf{A}^{-T}$
- ▶ Symmetric Matrices

$$\mathbf{A} = \mathbf{A}^T (a_{ij} = a_{ji}) \quad (5)$$

Special Cases of Matrix Multiplication: Inner and Outer Products

Given two vectors $\mathbf{a} \in \mathbb{R}^d$ and $\mathbf{b} \in \mathbb{R}^d$

- ▶ **Inner product:** $\mathbf{a}^T \mathbf{b} = \mathbf{a} \cdot \mathbf{b} = \langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^d a_i b_i.$
- ▶ **Outer product:** $\mathbf{a} \mathbf{b}^T = \mathbf{a} \otimes \mathbf{b} = [a_i b_j] \in \mathbb{R}^{d \times d}$

Special Cases of Matrix Multiplication: Inner and Outer Products

Given two vectors $\mathbf{a} \in \mathbb{R}^d$ and $\mathbf{b} \in \mathbb{R}^d$

- ▶ **Inner product:** $\mathbf{a}^T \mathbf{b} = \mathbf{a} \cdot \mathbf{b} = \langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^d a_i b_i.$
- ▶ **Outer product:** $\mathbf{a} \mathbf{b}^T = \mathbf{a} \otimes \mathbf{b} = [a_i b_j] \in \mathbb{R}^{d \times d}$

Given two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{m \times n}$

- ▶ **Inner product:** $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i,j} a_{ij} b_{ij}.$
- ▶ **Outer product:** $\mathbf{A} \otimes \mathbf{B} = [a_{ij} b_{kl}] \in \mathbb{R}^{m \times n \times m \times n}$

Special Cases of Matrix Multiplication: Inner and Outer Products

Given two vectors $\mathbf{a} \in \mathbb{R}^d$ and $\mathbf{b} \in \mathbb{R}^d$

- ▶ **Inner product:** $\mathbf{a}^T \mathbf{b} = \mathbf{a} \cdot \mathbf{b} = \langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^d a_i b_i$.
- ▶ **Outer product:** $\mathbf{a} \mathbf{b}^T = \mathbf{a} \otimes \mathbf{b} = [a_i b_j] \in \mathbb{R}^{d \times d}$

Given two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{m \times n}$

- ▶ **Inner product:** $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i,j} a_{ij} b_{ij}$.
- ▶ **Outer product:** $\mathbf{A} \otimes \mathbf{B} = [a_{ij} b_{kl}] \in \mathbb{R}^{m \times n \times m \times n}$

Kronecker product: $\mathbf{A} \otimes \mathbf{B} = [a_{ij} b_{kl}] \in \mathbb{R}^{m^2 \times n^2}$.

Some Properties of Matrix Multiplication

- ▶ Even if conformable, \mathbf{AB} does not necessarily equal \mathbf{BA} (i.e., matrix multiplication is not commutative)
- ▶ Matrix multiplication can be extended beyond two matrices
- ▶ matrix multiplication is associative, i.e., $\mathbf{A(BC)} = (\mathbf{AB})\mathbf{C}$
- ▶ Multiplication and transposition:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (6)$$

Orthogonal and Orthonormal

- ▶ If $\langle \mathbf{u}, \mathbf{v} \rangle = 0$, $\|\mathbf{u}\|_2 \neq 0$, $\|\mathbf{v}\|_2 \neq 0$, \mathbf{u} and \mathbf{v} are **orthogonal**.
- ▶ If $\langle \mathbf{u}, \mathbf{v} \rangle = 0$, $\|\mathbf{u}\|_2 = 1$, $\|\mathbf{v}\|_2 = 1$, \mathbf{u} and \mathbf{v} are **orthonormal**.

Orthogonal Matrix

- ▶ If square \mathbf{A} is orthogonal, it is easy to find its inverse:

$$\mathbf{A}\mathbf{A}^T = \mathbf{A}^T\mathbf{A} = \mathbf{I} \quad \left(\text{i.e., } \mathbf{A}^{-1} = \mathbf{A}^T \right) \quad (7)$$

- ▶ Property:

$$\|\mathbf{A}\mathbf{v}\| = \|\mathbf{v}\| \quad (\text{does not change the magnitude of } \mathbf{v}) \quad (8)$$

Determinant of Matrix

- ▶ The determinant of a matrix \mathbf{A} is denoted by $|\mathbf{A}|$ (or $\det(\mathbf{A})$ or $\det \mathbf{A}$).
- ▶ Determinants exist **only for square matrices**.
- ▶ 2×2

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \quad \det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{21}a_{12}$$

- ▶ 3×3

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

- ▶ $n \times n$

$$\det(\mathbf{A}) = \sum_{j=1}^m (-1)^{j+k} a_{jk} \det(\mathbf{A}_{jk}), \text{ for any } k : 1 \leq k \leq m$$

Determinant of Matrix

$$\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$$

$$\det(\mathbf{A} + \mathbf{B}) \neq \det(\mathbf{A}) + \det(\mathbf{B})$$

- ▶ diagonal matrix:

$$\text{If } \mathbf{A} = \begin{bmatrix} a_{11} & 0 & \cdot & 0 \\ 0 & a_{22} & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & a_{nn} \end{bmatrix}, \text{ then } \det(\mathbf{A}) = \prod_{i=1}^n a_{ii}$$

Inverse of Matrix

- ▶ The inverse of a matrix \mathbf{A} is commonly denoted by \mathbf{A}^{-1} or $\text{inv } \mathbf{A}$.
- ▶ The inverse of an $n \times n$ matrix \mathbf{A} is the matrix \mathbf{A}^{-1} such that $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I} = \mathbf{A}^{-1}\mathbf{A}$
- ▶ The matrix inverse is analogous to a scalar reciprocal
- ▶ A matrix which has an inverse is called **nonsingular**
- ▶ For some $n \times n$ matrix \mathbf{A} , an inverse matrix \mathbf{A}^{-1} may not exist.
- ▶ A matrix which does not have an inverse is **singular**.
- ▶ An inverse of $n \times n$ matrix \mathbf{A} exists if $|\mathbf{A}| \neq 0$.

Inverse of Matrix

- ▶ The inverse \mathbf{A}^{-1} of a matrix \mathbf{A} has the property: $\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$
- ▶ \mathbf{A}^{-1} exists if only if $\det(\mathbf{A}) \neq 0$
- ▶ Terminology
 - ▶ **Singular matrix:** \mathbf{A}^{-1} does not exist
 - ▶ **Ill-conditioned matrix:** \mathbf{A} is close to being singular

Property:

$$\begin{aligned}(\mathbf{AB})^{-1} &= \mathbf{B}^{-1}\mathbf{A}^{-1} \\ (\mathbf{A}^T)^{-1} &= (\mathbf{A}^{-1})^T \\ (\mathbf{A}^{-1})^{-1} &= \mathbf{A}\end{aligned}$$

- ▶ For diagonal matrices $\mathbf{D}^{-1} = \text{diag}\{d_1^{-1}, \dots, d_n^{-1}\}$
- ▶ For orthogonal matrices $\mathbf{A}^{-1} = \mathbf{A}^T$

Pseudo-inverse

- ▶ The pseudo-inverse \mathbf{A}^+ of a matrix \mathbf{A} (could be non-square, e.g., $m \times n$) is given by:

$$\mathbf{A}^+ = \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \mathbf{A}^T$$

- ▶ It can be shown that:

$$\mathbf{A}^+ \mathbf{A} = \mathbf{I} \quad (\text{provided that } \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \text{ exists})$$

Examples

- Fully-Connected Layer:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Weight}} \underbrace{\mathbf{X}}_{\text{Signals}} . \quad (9)$$

Examples

- Fully-Connected Layer:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Weight Signals}} \underbrace{\mathbf{X}}_{\text{Weight Signals}} . \quad (9)$$

- Graph Convolution:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Adj./Lap. Node Attr.}} \underbrace{\mathbf{X}}_{\text{Adj./Lap. Node Attr.}} . \quad (10)$$

Examples

- ▶ Fully-Connected Layer:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Weight Signals}} \underbrace{\mathbf{X}}_{\text{Weight Signals}} . \quad (9)$$

- ▶ Graph Convolution:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Adj./Lap. Node Attr.}} \underbrace{\mathbf{X}}_{\text{Adj./Lap. Node Attr.}} . \quad (10)$$

- ▶ Fast Fourier Transform:

$$\mathbf{B} = \underbrace{\mathbf{A}}_{\text{Basis Signals}} \underbrace{\mathbf{X}}_{\text{Basis Signals}} . \quad (11)$$

**Revisit above vector and matrix operations
from a viewpoint of linear systems.**

Key Tasks of Linear Systems

$$\underbrace{\mathbf{b}}_{\text{Output}} = \underbrace{\mathbf{A}}_{\text{System}} \underbrace{\mathbf{x}}_{\text{Input}} \quad (12)$$

- ▶ **Inverse Problem:** Both \mathbf{b} and \mathbf{A} are known,
 - ▶ Solve/Approximate the linear equation: $\mathbf{b} = \mathbf{A}\mathbf{x}$ or $\min_{\mathbf{x}} d(\mathbf{b}, \mathbf{A}\mathbf{x})$.

Key Tasks of Linear Systems

$$\underbrace{\mathbf{b}}_{\text{Output}} = \underbrace{\mathbf{A}}_{\text{System}} \underbrace{\mathbf{x}}_{\text{Input}} \quad (12)$$

- ▶ **Inverse Problem:** Both \mathbf{b} and \mathbf{A} are known,
 - ▶ Solve/Approximate the linear equation: $\mathbf{b} = \mathbf{A}\mathbf{x}$ or $\min_{\mathbf{x}} d(\mathbf{b}, \mathbf{A}\mathbf{x})$.
- ▶ **Modeling:** Given sets of \mathbf{b} 's and \mathbf{x} 's, denoted as \mathbf{B} and \mathbf{X} ,
 - ▶ Solve/Approximate the linear equation: $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}} d(\mathbf{B}, \mathbf{A}\mathbf{X})$.

Key Tasks of Linear Systems

$$\underbrace{\mathbf{b}}_{\text{Output}} = \underbrace{\mathbf{A}}_{\text{System}} \underbrace{\mathbf{x}}_{\text{Input}} \quad (12)$$

- ▶ **Inverse Problem:** Both \mathbf{b} and \mathbf{A} are known,
 - ▶ Solve/Approximate the linear equation: $\mathbf{b} = \mathbf{A}\mathbf{x}$ or $\min_{\mathbf{x}} d(\mathbf{b}, \mathbf{A}\mathbf{x})$.
- ▶ **Modeling:** Given sets of \mathbf{b} 's and \mathbf{x} 's, denoted as \mathbf{B} and \mathbf{X} ,
 - ▶ Solve/Approximate the linear equation: $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}} d(\mathbf{B}, \mathbf{A}\mathbf{X})$.
- ▶ **Factorization:** Given \mathbf{B} ,
 - ▶ Solve/Approximate the decomposition/factorization problem: $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}, \mathbf{X}} d(\mathbf{B}, \mathbf{A}\mathbf{X})$.

Key Tasks of Linear Systems

$$\underbrace{\mathbf{b}}_{\text{Output}} = \underbrace{\mathbf{A}}_{\text{System}} \underbrace{\mathbf{x}}_{\text{Input}} \quad (12)$$

- ▶ **Inverse Problem:** Both \mathbf{b} and \mathbf{A} are known,
 - ▶ Solve/Approximate the linear equation: $\mathbf{b} = \mathbf{A}\mathbf{x}$ or $\min_{\mathbf{x}} d(\mathbf{b}, \mathbf{A}\mathbf{x})$.
- ▶ **Modeling:** Given sets of \mathbf{b} 's and \mathbf{x} 's, denoted as \mathbf{B} and \mathbf{X} ,
 - ▶ Solve/Approximate the linear equation: $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}} d(\mathbf{B}, \mathbf{A}\mathbf{X})$.
- ▶ **Factorization:** Given \mathbf{B} ,
 - ▶ Solve/Approximate the decomposition/factorization problem: $\mathbf{B} = \mathbf{A}\mathbf{X}$ or $\min_{\mathbf{A}, \mathbf{X}} d(\mathbf{B}, \mathbf{A}\mathbf{X})$.

Many ML models, algorithms, and applications fall into these paradigms.

The main technical content of this course

Orthogonal Vectors and Linear Independence

- ▶ For $\mathbf{A} \in \mathbb{C}^{M \times N}$, its Hermitian transport (adjoint) is denoted as $\mathbf{A}^H \in \mathbb{C}^{N \times M}$.
- ▶ For $\mathbf{A} \in \mathbb{R}^{M \times N}$, its transport is $\mathbf{A}^T \in \mathbb{R}^{N \times M}$.

Orthogonal Vectors and Linear Independence

- ▶ For $\mathbf{A} \in \mathbb{C}^{M \times N}$, its Hermitian transport (adjoint) is denoted as $\mathbf{A}^H \in \mathbb{C}^{N \times M}$.
- ▶ For $\mathbf{A} \in \mathbb{R}^{M \times N}$, its transport is $\mathbf{A}^T \in \mathbb{R}^{N \times M}$.
- ▶ Inner product:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^H \mathbf{y} = \sum_{n=1}^N \bar{x}_n y_n, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^N \quad (13)$$

- ▶ **Explanation:** The projection of \mathbf{y} along the direction of \mathbf{x} .
- ▶ \mathbf{x} and \mathbf{y} are orthogonal to each other, if $\mathbf{x}^H \mathbf{y} = 0$.
- ▶ A **set** of *nonzero* vectors is orthogonal if its vectors are pairwise orthogonal.
- ▶ Orthonormal set of vectors: orthogonal + unit norm
- ▶ **Theorem.** Vectors in an orthogonal set are **linearly-independent**.

Orthogonal Vectors and Linear Independence

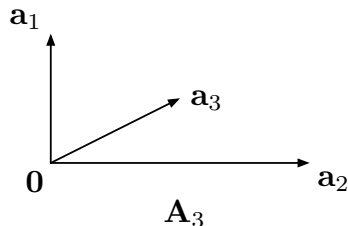
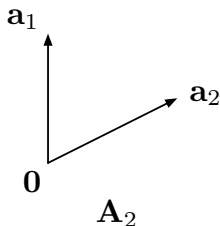
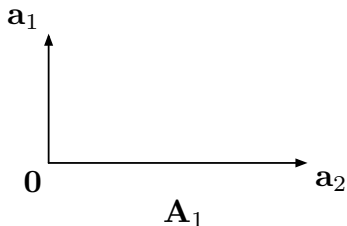
Definition (Linear Independence)

$[\mathbf{a}_1, \dots, \mathbf{a}_N]$ are linearly-independent $\Leftrightarrow \sum_{i=1}^N \mathbf{a}_i x_i = \mathbf{0}$ iff $\mathbf{x} = [x_i] = \mathbf{0}$.

Orthogonal Vectors and Linear Independence

Definition (Linear Independence)

$[\mathbf{a}_1, \dots, \mathbf{a}_N]$ are linearly-independent $\Leftrightarrow \sum_{i=1}^N \mathbf{a}_i x_i = \mathbf{0}$ iff $\mathbf{x} = [x_i] = \mathbf{0}$.

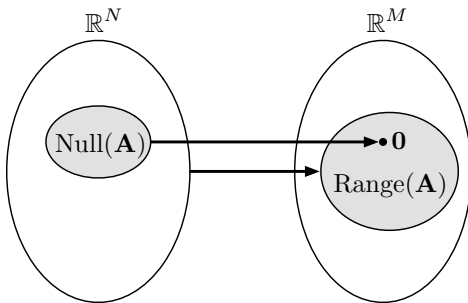


Range and Null Space

- ▶ $\text{Range}(\mathbf{A}) =$ The column space of $\mathbf{A} = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_N\} = \{\mathbf{A}\mathbf{x} | \mathbf{x} \in \mathbb{R}^N\}$.
- ▶ The null space of \mathbf{A} is $\text{Null}(\mathbf{A}) = \{\mathbf{x} | \mathbf{A}\mathbf{x} = \mathbf{0}\}$.
- ▶ Column rank: dimension of column space ($\leq N$).

Range and Null Space

- ▶ $\text{Range}(\mathbf{A}) =$ The column space of $\mathbf{A} = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_N\} = \{\mathbf{A}\mathbf{x} | \mathbf{x} \in \mathbb{R}^N\}$.
- ▶ The null space of \mathbf{A} is $\text{Null}(\mathbf{A}) = \{\mathbf{x} | \mathbf{A}\mathbf{x} = \mathbf{0}\}$.
- ▶ Column rank: dimension of column space ($\leq N$).



Rank of Matrix

- ▶ $\text{rank}(\mathbf{A})$ (the rank of a m -by- n matrix \mathbf{A}) is
 - = The maximal number of linearly independent columns
 - = The maximal number of linearly independent rows

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \text{Rank=?} \quad \begin{pmatrix} 2 & 1 \\ 4 & 2 \end{pmatrix} \text{Rank=?}$$

- ▶ If \mathbf{A} is n by m , then $\text{rank}(\mathbf{A}) \leq \min(m, n)$
- ▶ If $n = \text{rank}(\mathbf{A})$, then \mathbf{A} has full row rank
- ▶ If $m = \text{rank}(\mathbf{A})$, then \mathbf{A} has full column rank

Full Column Rank and Linear Independence

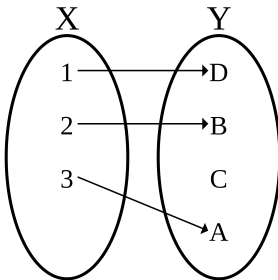
Theorem

\mathbf{A} is of full column rank ($\text{Rank}(\mathbf{A}) = N$).

$\Leftrightarrow \mathbf{A}$ is injective.

$\Leftrightarrow [\mathbf{a}_1, \dots, \mathbf{a}_N]$ is linearly-independent.

$\Leftrightarrow \text{Null}(\mathbf{A}) = \{\mathbf{0}\}$.



Linear (Vector) Space and Linear Independence

- ▶ Linear Vector Space: $\mathcal{X} \subset \mathbb{R}^N$.
- ▶ Computational closure:
 - ▶ Vector addition: $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$.
 - ▶ Scalar multiplication: $\mathcal{X} \times \mathbb{R} \mapsto \mathcal{X}$.

Linear (Vector) Space and Linear Independence

- ▶ Linear Vector Space: $\mathcal{X} \subset \mathbb{R}^N$.
- ▶ Computational closure:
 - ▶ Vector addition: $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$.
 - ▶ Scalar multiplication: $\mathcal{X} \times \mathbb{R} \mapsto \mathcal{X}$.
- ▶ The basis of a linear vector space \mathcal{X} : $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_D\}$
 - ▶ The \mathbf{b} 's are linear independent.
 - ▶ For each $\mathbf{x} \in \mathcal{X}$, $\exists \mathbf{a} \in \mathbb{R}^D$, $\mathbf{x} = \mathbf{B}\mathbf{a} = \sum_{i=1}^D \mathbf{b}_i a_i$.

Linear (Vector) Space and Linear Independence

- ▶ Linear Vector Space: $\mathcal{X} \subset \mathbb{R}^N$.
- ▶ Computational closure:
 - ▶ Vector addition: $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$.
 - ▶ Scalar multiplication: $\mathcal{X} \times \mathbb{R} \mapsto \mathcal{X}$.
- ▶ The basis of a linear vector space \mathcal{X} : $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_D\}$
 - ▶ The \mathbf{b} 's are linear independent.
 - ▶ For each $\mathbf{x} \in \mathcal{X}$, $\exists \mathbf{a} \in \mathbb{R}^D$, $\mathbf{x} = \mathbf{B}\mathbf{a} = \sum_{i=1}^D \mathbf{b}_i a_i$.
- ▶ **Question:** Can $D > N$?

Linear (Vector) Space and Linear Independence

- ▶ Linear Vector Space: $\mathcal{X} \subset \mathbb{R}^N$.
- ▶ Computational closure:
 - ▶ Vector addition: $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$.
 - ▶ Scalar multiplication: $\mathcal{X} \times \mathbb{R} \mapsto \mathcal{X}$.
- ▶ The basis of a linear vector space \mathcal{X} : $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_D\}$
 - ▶ The \mathbf{b} 's are linear independent.
 - ▶ For each $\mathbf{x} \in \mathcal{X}$, $\exists \mathbf{a} \in \mathbb{R}^D$, $\mathbf{x} = \mathbf{B}\mathbf{a} = \sum_{i=1}^D \mathbf{b}_i a_i$.
- ▶ **Question:** Can $D > N$? **No.**
- ▶ $D = \dim(\mathcal{X})$, and $\mathcal{X} = \text{span}(\mathbf{B})$.

Linear (Vector) Space and Linear Independence

- ▶ Linear Vector Space: $\mathcal{X} \subset \mathbb{R}^N$.
- ▶ Computational closure:
 - ▶ Vector addition: $\mathcal{X} \times \mathcal{X} \mapsto \mathcal{X}$.
 - ▶ Scalar multiplication: $\mathcal{X} \times \mathbb{R} \mapsto \mathcal{X}$.
- ▶ The basis of a linear vector space \mathcal{X} : $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_D\}$
 - ▶ The \mathbf{b} 's are linear independent.
 - ▶ For each $\mathbf{x} \in \mathcal{X}$, $\exists \mathbf{a} \in \mathbb{R}^D$, $\mathbf{x} = \mathbf{B}\mathbf{a} = \sum_{i=1}^D \mathbf{b}_i a_i$.
- ▶ **Question:** Can $D > N$? **No.**
- ▶ $D = \dim(\mathcal{X})$, and $\mathcal{X} = \text{span}(\mathbf{B})$.

Could you enumerate some typical linear spaces and their basis?

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$, where $\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ is an orthonormal set in \mathbb{R}^M .
- ▶ For a vector $\mathbf{v} \in \mathbb{R}^M$, its residual w.r.t. the set is

$$\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i. \quad (14)$$

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$, where $\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ is an orthonormal set in \mathbb{R}^M .
- ▶ For a vector $\mathbf{v} \in \mathbb{R}^M$, its residual w.r.t. the set is

$$\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i. \quad (14)$$

- ▶ Obviously, $\langle \mathbf{r}, \mathbf{q}_i \rangle = 0, \forall i = 1, \dots, N$. (Derive it)

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$, where $\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ is an orthonormal set in \mathbb{R}^M .
- ▶ For a vector $\mathbf{v} \in \mathbb{R}^M$, its residual w.r.t. the set is

$$\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i. \quad (14)$$

- ▶ Obviously, $\langle \mathbf{r}, \mathbf{q}_i \rangle = 0, \forall i = 1, \dots, N$. (Derive it)
- ▶ The decomposition of \mathbf{v} :

$$\mathbf{v} = \underbrace{\sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i}_{\in \mathcal{S}} + \underbrace{\mathbf{r}}_{\in \mathcal{S}^\perp}, \quad (15)$$

where $\mathcal{S} \oplus \mathcal{S}^\perp = \mathbb{R}^M$.

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$, where $\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ is an orthonormal set in \mathbb{R}^M .
- ▶ For a vector $\mathbf{v} \in \mathbb{R}^M$, its residual w.r.t. the set is

$$\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i. \quad (14)$$

- ▶ Obviously, $\langle \mathbf{r}, \mathbf{q}_i \rangle = 0, \forall i = 1, \dots, N$. (Derive it)
- ▶ The decomposition of \mathbf{v} :

$$\mathbf{v} = \underbrace{\sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i}_{\in \mathcal{S}} + \underbrace{\mathbf{r}}_{\in \mathcal{S}^\perp}, \quad (15)$$

where $\mathcal{S} \oplus \mathcal{S}^\perp = \mathbb{R}^M$.

- ▶ **Question:** $\forall \mathbf{x} \in \mathcal{S}$ and $\mathbf{y} \in \mathcal{S}^\perp$, is $\langle \mathbf{x}, \mathbf{y} \rangle \equiv 0$?

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$, where $\{\mathbf{q}_1, \dots, \mathbf{q}_N\}$ is an orthonormal set in \mathbb{R}^M .
- ▶ For a vector $\mathbf{v} \in \mathbb{R}^M$, its residual w.r.t. the set is

$$\mathbf{r} = \mathbf{v} - \sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i. \quad (14)$$

- ▶ Obviously, $\langle \mathbf{r}, \mathbf{q}_i \rangle = 0, \forall i = 1, \dots, N$. (Derive it)
- ▶ The decomposition of \mathbf{v} :

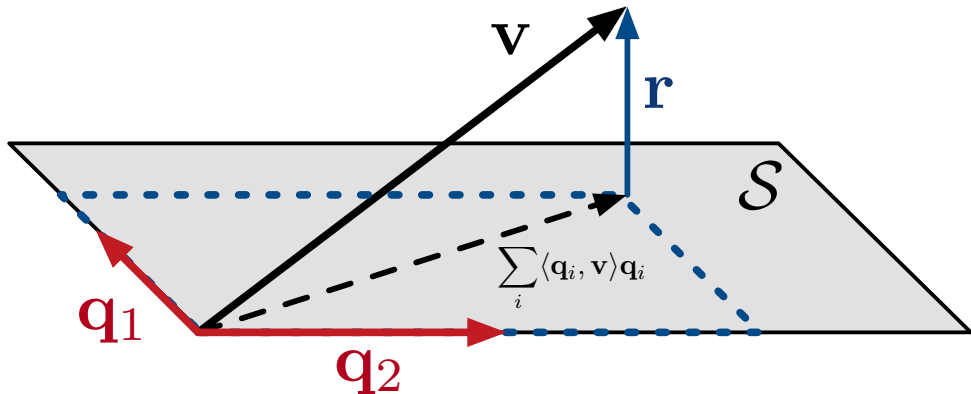
$$\mathbf{v} = \underbrace{\sum_{i=1}^N \langle \mathbf{q}_i, \mathbf{v} \rangle \mathbf{q}_i}_{\in \mathcal{S}} + \underbrace{\mathbf{r}}_{\in \mathcal{S}^\perp}, \quad (15)$$

where $\mathcal{S} \oplus \mathcal{S}^\perp = \mathbb{R}^M$.

- ▶ **Question:** $\forall \mathbf{x} \in \mathcal{S}$ and $\mathbf{y} \in \mathcal{S}^\perp$, is $\langle \mathbf{x}, \mathbf{y} \rangle \equiv 0$? **Yes.**

Components of a Vector

- ▶ Let $\mathcal{S} = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_2\}$ and $\mathbb{R}^M = \mathbb{R}^3$.



Recall The Linear System $\mathbf{y} = \mathbf{A}\mathbf{x}$

► $\mathbf{A} \in \mathbb{R}^{M \times M}, \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^M$

Recall The Linear System $\mathbf{y} = \mathbf{A}\mathbf{x}$

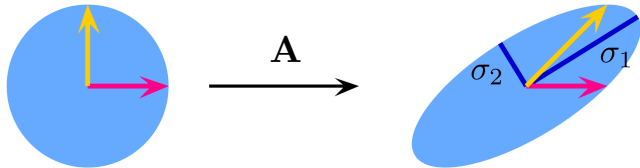
- ▶ $\mathbf{A} \in \mathbb{R}^{M \times M}$, $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^M$
- ▶ $\mathcal{Y} = \text{Range}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^M\}$ The image of the domain \mathcal{X} under the linear transform \mathbf{A} .

Recall The Linear System $\mathbf{y} = \mathbf{A}\mathbf{x}$

- ▶ $\mathbf{A} \in \mathbb{R}^{M \times M}$, $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^M$
- ▶ $\mathcal{Y} = \text{Range}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^M\}$ The image of the domain \mathcal{X} under the linear transform \mathbf{A} .
- ▶ When $\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^M \mid \|\mathbf{x}\|_2 = 1\}$ (a unit sphere \mathcal{S}^{M-1}), \mathcal{Y} is always a hyperellipse.

Recall The Linear System $\mathbf{y} = \mathbf{A}\mathbf{x}$

- ▶ $\mathbf{A} \in \mathbb{R}^{M \times M}$, $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^M$
- ▶ $\mathcal{Y} = \text{Range}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^M\}$ The image of the domain \mathcal{X} under the linear transform \mathbf{A} .
- ▶ When $\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^M \mid \|\mathbf{x}\|_2 = 1\}$ (a unit sphere \mathcal{S}^{M-1}), \mathcal{Y} is always a hyperellipse.
- ▶ A hyperellipse $\mathcal{Y} = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathcal{S}^{M-1}\}$, characterized by $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{R}^M$ orthonormal directions, and $\sigma_1, \dots, \sigma_M$ the corresponding length of the semi-axes.



Derivative of a Function

- ▶ $\lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$ is called the derivative of f at a .
- ▶ We write the derivative of f with respect to x is

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

- ▶ There are many ways to write the derivative of $y = f(x)$.
e.g. define the slope of the curve $y=f(x)$ at the point x .

Some Important Rules of (Partial) Derivatives

- ▶ Scalar multiplication: $\partial_x[af(x)] = a [\partial_x f(x)]$
- ▶ Polynomials: $\partial_x [x^k] = kx^{k-1}$
- ▶ Function addition: $\partial_x[f(x) + g(x)] = [\partial_x f(x)] + [\partial_x g(x)]$
- ▶ Function multiplication: $\partial_x[f(x)g(x)] = f(x) [\partial_x g(x)] + [\partial_x f(x)] g(x)$
- ▶ Function division: $\partial_x \left[\frac{f(x)}{g(x)} \right] = \frac{[\partial_x f(x)]g(x) - f(x)[\partial_x g(x)]}{[g(x)]^2}$
- ▶ Function composition: $\partial_x[f(g(x))] = [\partial_x f] (g(x)) [\partial_x g(x)]$
- ▶ Exponentiation: $\partial_x [e^x] = e^x$ and $\partial_x [a^x] = \log(a)e^x$
- ▶ Logarithms: $\partial_x[\log x] = \frac{1}{x}$

Definitions of Gradient

- ▶ Matrix-calculus Scalar-by-matrix
- ▶ Suppose that $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a function that takes as input a matrix A of size $m \times n$ and returns a real value. Then the **gradient** of f (with respect to $\mathbf{A} \in \mathbb{R}^{m \times n}$) is the matrix of

$$\nabla_{\mathbf{A}} f(\mathbf{A}) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(\mathbf{A})}{\partial A_{11}} & \frac{\partial f(\mathbf{A})}{\partial A_{12}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial A_{1n}} \\ \frac{\partial f(\mathbf{A})}{\partial A_{21}} & \frac{\partial f(\mathbf{A})}{\partial A_{22}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{A})}{\partial A_{m1}} & \frac{\partial f(\mathbf{A})}{\partial A_{m2}} & \cdots & \frac{\partial f(\mathbf{A})}{\partial A_{mn}} \end{bmatrix}$$

- ▶ In principle, gradients are a natural extension of partial derivatives to functions of multiple variables.

Definitions of Gradient

- ▶ $f: \mathcal{X} \subset \mathbb{R}^n \mapsto \mathbb{R}$
- ▶ Size of gradient is always the same as the size of variable

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \text{ if } \mathbf{x} \in \mathbb{R}^n \quad (16)$$

Definitions of Gradient

- ▶ $f: \mathcal{X} \subset \mathbb{R}^n \mapsto \mathbb{R}^m$
- ▶ Size of gradient is always the same as the size of variable

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}$$

Definitions of Gradient

- ▶ $f: \mathcal{X} \subset \mathbb{R}^n \mapsto \mathbb{R}^m$
- ▶ Size of gradient is always the same as the size of variable

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} = \mathbf{J}_f^T(\mathbf{x}) \quad (17)$$

Examples

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T$$

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \left(\mathbf{B} + \mathbf{B}^T \right) \mathbf{x}$$

Hessian Matrix

- Suppose that $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that takes a vector in \mathbb{R}^n and returns a real number. Then the Hessian matrix with respect to \mathbf{x} , written $\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or simply as \mathbf{H} is the $n \times n$ matrix of partial derivatives,

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix} \quad (18)$$

References

- ▶ More knowledge about matrix analysis will be introduced later
 - ▶ SVD
 - ▶ Eigen decomposition
- ▶ Matrix Cookbook

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Statistical Machine Learning: Connecting Statistics and Algebra

A typical ML scenario:

- ▶ $\mathbf{X} \in \mathbb{R}^{D \times N}$ are a set of samples.
- ▶ Each sample $\mathbf{x} \sim \mu_{\mathcal{X}}$ is a **random variable**.
- ▶ How to estimate $P_{\mathcal{X}}$ via a model \hat{p}_{θ} based on the data \mathbf{X} ?

Mean, (Co)variance, and Their Unbiased Estimation

Suppose that we observed a set of i.i.d. samples $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$, each $\mathbf{x}_i \sim P$

- Mean

$$\mu = \mathbb{E}_P[X] \tag{19}$$

- Unbiased estimation of mean: the average of samples

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \tag{20}$$

Mean, (Co)variance, and Their Unbiased Estimation

Suppose that we observed a set of i.i.d. samples $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$, each $\mathbf{x}_i \sim P$

- Mean

$$\mu = \mathbb{E}_P[X] \quad (19)$$

- Unbiased estimation of mean: the average of samples

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (20)$$

- Variance

$$\sigma^2 = \mathbb{V}_P[X] = \mathbb{E}_P[(X - \mu)^2] = \underbrace{\mathbb{E}_P[X^2] - \mathbb{E}_P^2[X]}_{\text{Derive it}} \quad (21)$$

- Unbiased estimation of variance:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2 \quad \text{or} \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})^2 \quad (22)$$

Properties of Mean and Variance

- ▶ $\mathbb{E}[aX] = a\mathbb{E}[X]$
- ▶ $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$
- ▶ $\mathbb{V}[X + a] = \mathbb{V}[X]$
- ▶ $\mathbb{V}[aX + bY] = a^2\mathbb{V}[X] + b^2\mathbb{V}[Y] + 2ab\text{Cov}(X, Y).$

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \quad (23)$$

- ▶ Linear combinations

$$\mathbb{V}\left[\sum_{i=1}^K X_i\right] = \sum_{i,j=1}^K \text{Cov}(X_i, X_j). \quad (24)$$

Method of Moments

- ▶ Suppose that we have a model $\theta \in \mathbb{R}^D$, which works to define a data distribution $P(X; \theta)$.
- ▶ The parameters can be determined by solving the equations corresponding to the top- D moments of the distribution

$$\mu_d = \mathbb{E}_P[X^d] = f_d(\theta), \quad d = 1, \dots, D. \quad (25)$$

- ▶ In practice, given a set of samples $\mathbf{X} = [x_i]$, we can estimate $\{\mu_d\}$ as

$$\hat{\mu}_d = \frac{1}{n} \sum_{i=1}^n x_i^d, \quad d = 1, \dots, D. \quad (26)$$

Ideally, solving the D equations in (25) provides us with a good estimation of θ .

Method of Moments

- ▶ Suppose that we have a model $\theta \in \mathbb{R}^D$, which works to define a data distribution $P(X; \theta)$.
- ▶ The parameters can be determined by solving the equations corresponding to the top- D moments of the distribution

$$\mu_d = \mathbb{E}_P[X^d] = f_d(\theta), \quad d = 1, \dots, D. \quad (25)$$

- ▶ In practice, given a set of samples $\mathbf{X} = [x_i]$, we can estimate $\{\mu_d\}$ as

$$\hat{\mu}_d = \frac{1}{n} \sum_{i=1}^n x_i^d, \quad d = 1, \dots, D. \quad (26)$$

Ideally, solving the D equations in (25) provides us with a good estimation of θ .

- ▶ **Could you enumerate its drawbacks?**

Why Does ML Require Lots of Data: Law of Large Numbers

Weak law (converge in probability)

- For $\{X_i\}_{i=1}^n$,

$$\lim_{n \rightarrow \infty} \bar{X}_n = \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n X_i}{n} \xrightarrow{P} \underbrace{\mathbb{E}[X]}_{\mu}. \quad (27)$$

Why Does ML Require Lots of Data: Law of Large Numbers

Weak law (converge in probability)

- For $\{X_i\}_{i=1}^n$,

$$\lim_{n \rightarrow \infty} \bar{X}_n = \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n X_i}{n} \xrightarrow{P} \underbrace{\mathbb{E}[X]}_{\mu}. \quad (27)$$

- Equivalent representation

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| < \epsilon) = 1, \quad \forall \epsilon > 0. \quad (28)$$

Why Does ML Require Lots of Data: Law of Large Numbers

Strong law (Kolmogorov's law, converge almost surely)

- For $\{X_i\}_{i=1}^n$,

$$\lim_{n \rightarrow \infty} \bar{X}_n \xrightarrow{a.s.} \mu \quad (29)$$

or equivalently,

$$P \left(\lim_{n \rightarrow \infty} \bar{X}_n = \mu \right) = 1. \quad (30)$$

Why Does ML Require Lots of Data: Law of Large Numbers

Strong law (Kolmogorov's law, converge almost surely)

- For $\{X_i\}_{i=1}^n$,

$$\lim_{n \rightarrow \infty} \bar{X}_n \xrightarrow{a.s.} \mu \quad (29)$$

or equivalently,

$$\mathbb{P} \left(\lim_{n \rightarrow \infty} \bar{X}_n = \mu \right) = 1. \quad (30)$$

Variance reduction: Suppose that $\mathbb{V}[X_i] = \sigma^2$ for $i = 1, 2, \dots$

$$\mathbb{V}[\bar{X}_n] = \mathbb{V} \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n^2} \mathbb{V} \left[\sum_{i=1}^n X_i \right] = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}. \quad (31)$$

Why Does ML Like Gaussian Distribution: Central Limit Theorem

Lindeberg-Lévy CLT.

- Suppose $\{X_i\}_{i=1}^n$ is a sequence of i.i.d. random variables with $\mathbb{E}[X_i] = \mu$ and $\mathbb{V}[X_i] = \sigma^2 < \infty$, then

$$\lim_{n \rightarrow \infty} \sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2) \quad (32)$$

where \xrightarrow{d} means converge in distribution.

ML: Frequentist Statistic Viewpoint

Given a set of samples $\mathbf{X} = [\mathbf{x}_i]$, we assume that the samples are sampled from a distribution $P(\mathbf{x}|\theta)$ (a model with parameter θ).

- **Principle:** Assume a **deterministic** model, learn the model via maximum likelihood estimation (**MLE**):

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} \underbrace{\log P(\mathbf{X}|\theta)}_{\text{Loglike}(\theta; \mathbf{X})} \\ &= \arg \max_{\theta} \sum_{i=1}^n \log P(\mathbf{x}_i|\theta) \quad (\text{i.i.d. Assumption})\end{aligned}\tag{33}$$

ML: Frequentist Statistic Viewpoint

Given a set of samples $\mathbf{X} = [\mathbf{x}_i]$, we assume that the samples are sampled from a distribution $P(\mathbf{x}|\theta)$ (a model with parameter θ).

- **Principle:** Assume a **deterministic** model, learn the model via maximum likelihood estimation (**MLE**):

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} \underbrace{\log P(\mathbf{X}|\theta)}_{\text{Loglike}(\theta; \mathbf{X})} \\ &= \arg \max_{\theta} \sum_{i=1}^n \log P(\mathbf{x}_i|\theta) \quad (\text{i.i.d. Assumption})\end{aligned}\tag{33}$$

- Deterministic model \rightarrow (pointwise) MLE \rightarrow optimization.

ML: Bayesian Statistic Viewpoint

- ▶ **Principle:** Assume a **probabilistic** model — the model θ yields a prior distribution.
- ▶ Bayes' theorem

$$\underbrace{P(\theta|\mathbf{X})}_{\text{Posterior}(\theta)} = \frac{P(\mathbf{X}|\theta)P(\theta)}{P(\mathbf{X})} \propto \underbrace{P(\mathbf{X}|\theta)}_{\text{Likelihood}(\theta)} \underbrace{P(\theta)}_{\text{Prior}(\theta)} . \quad (34)$$

- ▶ Maximum A Posterior (**MAP**) estimation

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|\mathbf{X}) = \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta) \quad (35)$$

ML: Bayesian Statistic Viewpoint

- ▶ **Principle:** Assume a **probabilistic** model — the model θ yields a prior distribution.
- ▶ Bayes' theorem

$$\underbrace{P(\theta|\mathbf{X})}_{\text{Posterior}(\theta)} = \frac{P(\mathbf{X}|\theta)P(\theta)}{P(\mathbf{X})} \propto \underbrace{P(\mathbf{X}|\theta)}_{\text{Likelihood}(\theta)} \underbrace{P(\theta)}_{\text{Prior}(\theta)}. \quad (34)$$

- ▶ Maximum A Posterior (**MAP**) estimation

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|\mathbf{X}) = \arg \max_{\theta} P(\mathbf{X}|\theta)P(\theta) \quad (35)$$

- ▶ Probabilistic model \rightarrow (distribution) MAP \rightarrow optimization (Variational Inference) or sampling (MCMC).
- ▶ The influence of prior decays with the increase of the number of samples.

Frequentist v.s. Bayesian

Frequentist

- ▶ **Pros:** more efficient in general, avoid the design of prior.
- ▶ **Cons:** non-robust to sparse data, cannot quantify the uncertainty of the estimation.

Bayesian

- ▶ **Pros:** prior makes it (relatively) robust to sparse data, quantify the uncertainty of the estimation (obtain the distribution of θ)
- ▶ **Cons:** require sophisticated design of prior, time-consuming in general.

In Summary

- ▶ We review more matrix operations and their properties
- ▶ A viewpoint of linear algebra is provided and connected with machine learning
- ▶ Some statistical concepts are provided and connected with matrix operations and machine learning

Next...

- ▶ Linear regression model
- ▶ Learning, evaluation, and some theoretical results.