

Review Session 1

Zhiyuan Peng, Dexu Kong, Wanda Li, Weida Wang

September 13, 2022

Aim: This note is to review some basic mathematical knowledge on linear algebra, calculus and probability. We hope it can assist you in your future coursework.

1 Linear Algebra

1.1 Inner Product and trace

Definition 1. (Inner product). A function $\langle \cdot, \cdot \rangle: \mathbb{V} \times \mathbb{V} \rightarrow \mathbb{F}$ is an inner product if it satisfies [1]:

- **Linearity:** $\langle \alpha \mathbf{v} + \beta \mathbf{w}, \mathbf{x} \rangle = \alpha \langle \mathbf{v}, \mathbf{x} \rangle + \beta \langle \mathbf{w}, \mathbf{x} \rangle$;
- **Conjugate symmetry:** $\langle \mathbf{v}, \mathbf{w} \rangle = \overline{\langle \mathbf{w}, \mathbf{v} \rangle}$;
- **Positive definiteness:** $\langle \mathbf{v}, \mathbf{v} \rangle \geq 0$, with the equality iff $\mathbf{v} = \mathbf{0}$.

The most common one is the canonical inner product on \mathbb{R}^n . It says for vectors $\mathbf{x} \triangleq [x_1, \dots, x_n]^T$ and $\mathbf{y} \triangleq [y_1, \dots, y_n]^T$, we have

$$\langle \mathbf{x}, \mathbf{y} \rangle \triangleq x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{i=1}^n x_i y_i = \mathbf{x}^T \mathbf{y}.$$

Example 1: (Orthogonal Vectors) Vector $\mathbf{x} \in \mathbb{R}^n$ is orthogonal to $\mathbf{y} \in \mathbb{R}^n$ when $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

Example 2: (Unit Vector) Vector $\mathbf{x} \in \mathbb{R}^n$ is of unit length when $\langle \mathbf{x}, \mathbf{x} \rangle = 1$.

Example 3: (Orthogonal Matrix) The matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is said to be **orthogonal** if

$$\mathbf{Q}\mathbf{Q}^T = \mathbf{Q}^T\mathbf{Q} = \mathbf{I}$$

which implies that each column of \mathbf{Q} has unit length and orthogonal to each other.

Definition 2. (Trace). For $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\text{trace}(\mathbf{M}) = \sum_{i=1}^n \mathbf{M}_{ii}$, where \mathbf{M}_{ii} is the diagonal terms of matrix \mathbf{M} .

Theorem 1. For any matrices \mathbf{A}, \mathbf{B} of compatible size,

$$\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA}).$$

1.2 Vector Norms

A norm on a vector space \mathbb{V} gives a way of measuring lengths of vectors. Formally:

Definition 3. (Vector norm). A norm on a real vector space \mathbb{V} is a function $\|\cdot\| : \mathbb{V} \rightarrow \mathbb{R}$ that is:

- **Nonnegatively homogeneous:** $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ for all vectors $\mathbf{x} \in \mathbb{V}$, scalars $\alpha \in \mathbb{R}$;
- **Positive definite:** $\|\mathbf{x}\| \geq 0$, and $\|\mathbf{x}\| = 0$ iff $\mathbf{x} = \mathbf{0}$;
- **Subadditive:** $\|\cdot\|$ satisfies the triangle inequality $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$, for all $\mathbf{x}, \mathbf{y} \in \mathbb{V}$.

One very important family of norms are the ℓ^p norms. If we take $\mathbb{V} = \mathbb{R}^n$, and $p \in [1, \infty)$, for vector $\mathbf{x} \triangleq [x_1, \dots, x_n]^T$, we have

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}. \quad (1)$$

The most frequent used one is the ℓ^2 norm or the "Euclidean norm",

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\mathbf{x}^T \mathbf{x}}$$

which coincides with our usually way of measuring lengths. Two other cases are of almost equal importance: $p = 1$, and $p \rightarrow \infty$. Setting $p = 1$ in (1), we obtain $\|\mathbf{x}\|_1 = \sum_i |x_i|$.

Finally, as p becomes larger, the expression in (1) accentuates the largest $|x_i|$ among \mathbf{x} entries. In another words, as $p \rightarrow \infty$, $\|\mathbf{x}\|_p \rightarrow \max_i |x_i|$. Thus, we can extend the definition of the ℓ^p norm to $p = \infty$ by defining

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

2 Calculus

2.1 Derivatives

Scalar b , vectors $\mathbf{x}, \mathbf{w}, \mathbf{y}$ and matrix \mathbf{A} , we have :

- $\frac{\partial(\mathbf{w}^T \mathbf{x} + b)}{\partial \mathbf{x}} = \mathbf{w}$
- $\frac{\partial(\mathbf{x}^T \mathbf{A} \mathbf{x} + b)}{\partial \mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x}$
- $\frac{\partial(\mathbf{x}^T \mathbf{A}^{-1} \mathbf{y})}{\partial \mathbf{A}} = -\mathbf{A}^{-T} \mathbf{x} \mathbf{y}^T \mathbf{A}^{-T}$

For more derivative calculation, please refer to the [Matrix Cookbook](#)[2].

3 Probability

3.1 Basic Properties

For events E_1 and E_2 , if they are disjoint, i.e. $E_1 \cap E_2 = \emptyset$, then $\mathbb{P}(E_1 \cup E_2) = \mathbb{P}(E_1) + \mathbb{P}(E_2)$

Definition 4. (Conditional probability) For events A and B , and $\mathbb{P}(A) > 0$,

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)}$$

We can define the conditional expectation as

$$\mathbb{E}[Y|X = x] \triangleq \sum_{y \in \mathcal{Y}} y \cdot p(Y = y|X = x)$$

Definition 5. (Covariance) For two random variables X and Y , the covariance is defined by

$$\text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

When the covariance of X and Y is 0, we call them **uncorrelated variables**.

Definition 6. (Independent) For two random variables, when the joint pdf can be written as the product of two RVs' pdf

$$f(x, y) = f_X(x) f_Y(y),$$

we call them **independent**.

Theorem 2. We have:

◦ (Multiplication Rule) For events A and B ,

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B|A) = \mathbb{P}(B)\mathbb{P}(A|B);$$

◦ (Total probability rule) B_1, B_2, \dots, B_k form a partition of Ω , $\forall i \neq j, B_i \cap B_j = \emptyset, \cup_{i=1}^k B_i = \Omega$, we have:

$$\mathbb{P}(A) = \sum_{i=1}^k \mathbb{P}(B_i)\mathbb{P}(A|B_i);$$

◦ (Bayes Rule)

$$\mathbb{P}(B_1|A) = \frac{\mathbb{P}(A \cap B_1)}{\mathbb{P}(A)} = \frac{\mathbb{P}(A|B_1)\mathbb{P}(B_1)}{\mathbb{P}(A)} = \frac{\mathbb{P}(A|B_1)\mathbb{P}(B_1)}{\sum_{i=1}^k \mathbb{P}(A|B_i)\mathbb{P}(B_i)}.$$

3.2 Gaussian Distribution

3.2.1 Normal Distribution

• If random variable $X \in \mathbb{R}$, $X \sim \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and $\sigma \in \mathbb{R}$, then the density function of it is:

$$p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

• $\mathbb{E}[X] = \mu$; $\text{var}(X) = \sigma^2$.

3.2.2 Multivariate Gaussian Distribution

• If random variable $\mathbf{X} \in \mathbb{R}^n$, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ is symmetric and positive semi-definite (PSD), then the density function of it is:

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

• $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$; $\text{cov}(\mathbf{X}) = \boldsymbol{\Sigma}$.

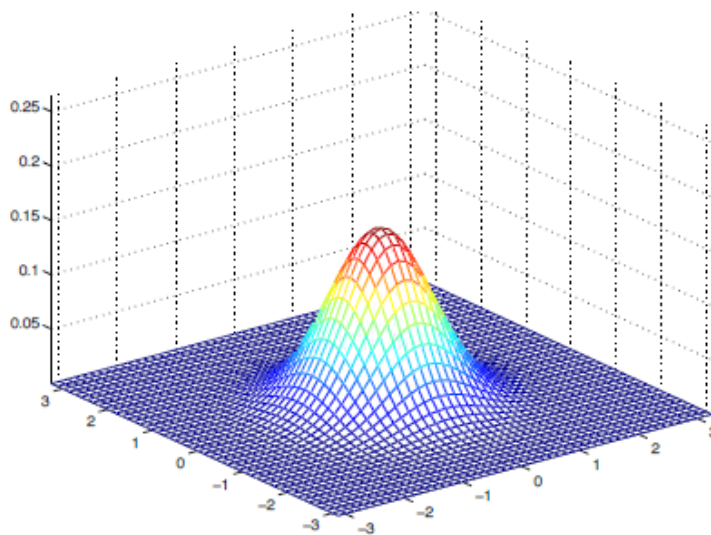


Figure 1: Multivariate Gaussian's p.d.f

References

- [1] Strang, Gilbert, et al. Introduction to linear algebra. Vol. 3. Wellesley, MA: Wellesley-Cambridge Press, 1993.
- [2] The Matrix Cookbook <http://matrixcookbook.com>