

MATH3332 Data Analytic Tools

Ye Moe

HKUST Fall 2022

Introduction

The purpose of this course is to introduce some crucial mathematical analysis tools for data analysis/machine learning.

According to *Pedro Domingos*,

$$\textit{Learning} = \textit{Representation} + \textit{Evaluation} + \textit{Optimization}$$

1. Representation

- How do we represent a learner? Which set should a learner be in? This set is called the hypothesis space of the learner. Some related tools are "space of functions".
- How do we represent the input? Potential tools include vectors, graphs, manifolds, ...

2. Evaluation

- How to pick the best learner from the hypothesis space? Needs calculus of "functions of functions" also known as functionals.
- How to represent the input effectively? Needs Linear Algebra, Graph Theory, Manifolds Calculus, Harmonic Analysis, ...

3. Optimization

- Numerical optimization solver - how to get the optimal solution numerically by a computer? Many of the resulting optimization is convex optimization and it is related to Convex Analysis.

So this course consists of some

- Basic functional analysis (calculus of functionals)
- Basic convex analysis
- Fourier analysis and Wavelet analysis (if time allowed)

Normed and Inner Product Space

2.1 Vector Spaces

Definition: A vector space over \mathbb{R} is a set \mathbb{V} together with two functions.

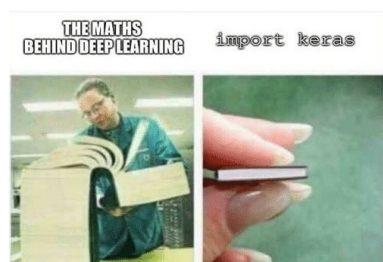
1. Vector addition: $+: (\mathbb{V}, \mathbb{V}) \rightarrow \mathbb{V}$
i.e. $\forall x, y \in \mathbb{V}, x + y \in \mathbb{V}$
2. Scalar multiplication: $\cdot: (\mathbb{R}, \mathbb{V}) \rightarrow \mathbb{V}$
i.e. $\forall \alpha \in \mathbb{R}, x \in \mathbb{V}, \alpha x \in \mathbb{V}$

These two functions should satisfy the following eight properties:

1. Associativity of addition: $x + (y + z) = (x + y) + z, \forall x, y, z \in \mathbb{V}$
2. Commutativity of addition: $x + y = y + x, \forall x, y \in \mathbb{V}$
3. Zero vector: \exists an element, denoted by 0 in \mathbb{V} s.t. $x + 0 = 0 + x = x, \forall x \in \mathbb{V}$
4. Negative vector: $\forall x \in \mathbb{V}, \exists$ an elements, denoted by $-x \in \mathbb{V}$ s.t. $x + (-x) = (-x) + x = 0$
5. $\forall x \in \mathbb{V}, 1 \cdot x = x$
6. $\forall x \in \mathbb{V}, \alpha, \beta \in \mathbb{R}, \alpha(\beta x) = (\alpha\beta)x$
7. $\forall x \in \mathbb{V}$ and $\alpha, \beta \in \mathbb{R}, (\alpha + \beta)x = \alpha x + \beta x$
8. $\forall x, y \in \mathbb{V}, \alpha(x + y) = \alpha x + \alpha y$

Remarks: We can define vector space over the complex domain \mathbb{C} , but since vector space over complex domain \mathbb{C} is used very rarely, we will only consider vector space in the real domain \mathbb{R} .

Some examples of vector space include $\mathbb{R}, \mathbb{R}^n, \mathbb{R}^{m \times n}, \mathbb{R}^{m \times n \times l}, C[a, b]$ and L_∞ .



Machine learning be like

Example: Prove that \mathbb{R}^n is a vector space.
 $\forall x, y \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$,

$$x + y = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix} \in \mathbb{R}^n$$

$$\alpha x = \alpha \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{bmatrix} \in \mathbb{R}^n$$

Since it is closed under both vector addition and scalar multiplication, \mathbb{R}^n is a vector space.

Example: Prove that $C[a, b]$ is a vector space.
 $\forall f, g \in C[a, b]$ and $\alpha \in \mathbb{R}$,

$$f(t) + g(t) = (f + g)(t) \in C[a, b], \forall t \in [a, b]$$

$$\alpha f(t) = (\alpha f)(t) \in C[a, b], \forall t \in [a, b]$$

Since it is closed under both vector addition and scalar multiplication, $C[a, b]$ is a vector space.

Remarks: $C[a, b]$ is referred to as a function space, since any vector in this vector space is a function. It might be a hypothesis space of a learner with one input and one output, i.e. Find a $f \in C[a, b]$ s.t. $f(x_i) \approx f(y_i)$ for all i .

Example: Prove that L_∞ is a vector space.

$$L_\infty = \left\{ \begin{bmatrix} a_1 \\ a_2 \\ \vdots \end{bmatrix} \mid \exists \text{ a finite number } c \text{ s.t. } |a_i| \leq c \text{ for any } i \right\}$$

$\forall a, b \in L_\infty$ and $\alpha \in \mathbb{R}$,

$$a + b = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} a_1 + b_1 \\ a_2 + b_2 \\ \vdots \end{bmatrix} \in L_\infty$$

$$\alpha a = \alpha \begin{bmatrix} a_1 \\ a_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} \alpha a_1 \\ \alpha a_2 \\ \vdots \end{bmatrix} \in L_\infty$$

Since it is closed under both vector addition and scalar multiplication, L_∞ is a vector space.

Remarks: This vector space can be used to model stock prices with a very fine time resolution.

Example: Consider the set of all strings.

$$'I' + 'am' \neq 'am' + 'I'$$

The set of all strings violates the commutative properties of a vector space, therefore it isn't a vector space. Hence, we cannot use vector space to model text data in this naïve way.

How do we "vectorize" the text data?
This is a fundamental question in text data analysis.

2.2 Normed and Banach Space

In order to do calculus on vector spaces, we need to define 'distance/closeness' between vectors.

Let \mathbb{V} be a vector space. Let $x, y \in \mathbb{V}$. Then,

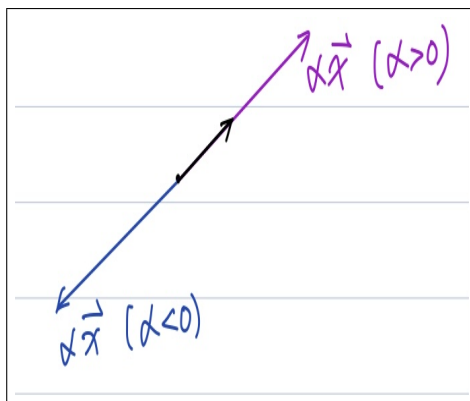
$$\text{distance}(x, y) = \text{distance}(x - y, y - y) = \text{distance}(x - y, 0) = \text{length of } x - y$$

Remarks: Distance should be shift invariant.

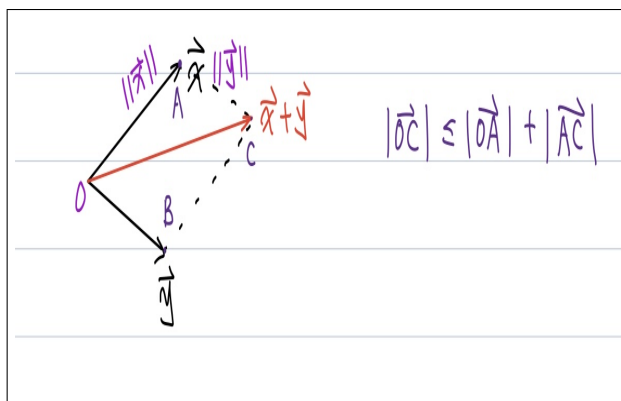
To define distance, we only need to define the length of vectors. Let $x \in \mathbb{V}$.

Denote $\|x\|$ be the length of x . Then $\|x\|$ should satisfy:

1. $\|x\| \geq 0$ (the length should be non-negative)
 Moreover, $\|x\| = 0 \iff x = 0$ (only zero vector has a zero length)
2. $\|\alpha x\| = |\alpha| \|x\|$
 (length of a scaling of a vector is a scaling of the length of the vector)



3. $\|x + y\| \leq \|x\| + \|y\|$ (also known as triangle inequality)
 (length of direct path should be smaller than the length of indirect path)



Definition: Let \mathbb{V} be a vector space. A norm on \mathbb{V} is a function $\|\cdot\| : \mathbb{V} \rightarrow \mathbb{R}$ such that:

1. $\|x\| \geq 0 \ \forall x \in \mathbb{V}$ and $\|x\| = 0 \iff x = 0$
2. $\|\alpha x\| = |\alpha| \|x\|, \ \forall \alpha \in \mathbb{R}, x \in \mathbb{V}$
3. $\|x + y\| \leq \|x\| + \|y\|, \ \forall x, y \in \mathbb{V}$

Example: \mathbb{R} is a vector space over \mathbb{R} .
Let $\|x\| = |x| \ \forall x \in \mathbb{R}$. Then it is a norm on \mathbb{R} .

Example: \mathbb{R}^n is a vector space over \mathbb{R} .
There are many norms on \mathbb{R}^n .

- 2-norm: (Euclidean Norm)
 $\|x\|_2 = (\sum_{i=1}^n x_i^2)^{\frac{1}{2}}$

Question: Prove that $\|\cdot\|_2$ is indeed a norm for \mathbb{R}^n .
 $\forall x, y \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$,

$$\|x\|_2 = (\sum_{i=1}^n x_i^2)^{\frac{1}{2}} \geq 0$$

$$\|x\|_2 = 0 \iff \sum_{i=1}^n x_i^2 = 0 \iff x_i^2 = 0, \ i = 1, \dots, n$$

$$\iff x_i = 0, \ i = 1, \dots, n \iff x = 0$$

$$\|\alpha x\|_2 = (\sum_{i=1}^n (\alpha x_i)^2)^{\frac{1}{2}} = (\alpha^2 \sum_{i=1}^n x_i^2)^{\frac{1}{2}} = |\alpha| (\sum_{i=1}^n x_i^2)^{\frac{1}{2}} = |\alpha| \|x\|_2$$

$$\|x + y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 + 2\langle x, y \rangle$$

$$\leq \|x\|_2^2 + \|y\|_2^2 + 2\|x\|_2\|y\|_2 \text{ (By Cauchy-Schwartz inequality)}$$

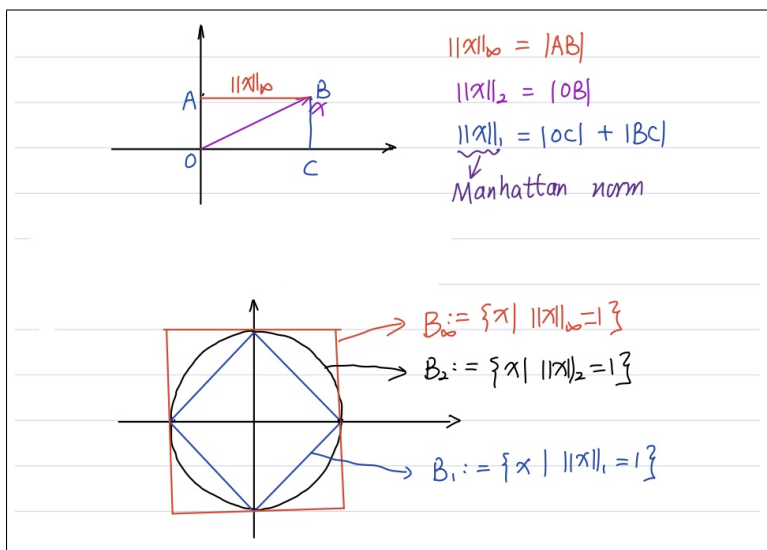
$$= (\|x\|_2 + \|y\|_2)^2$$

$$\|x + y\|_2 \leq \|x\|_2 + \|y\|_2$$

- 1-norm:
 $\|x\|_1 = \sum_{i=1}^n |x_i|$
- ∞ -norm:
 $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$
- p-norm:
 $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}$

Fact: $\|x\|_p$ is a norm on $\mathbb{R}^n \iff p \geq 1$

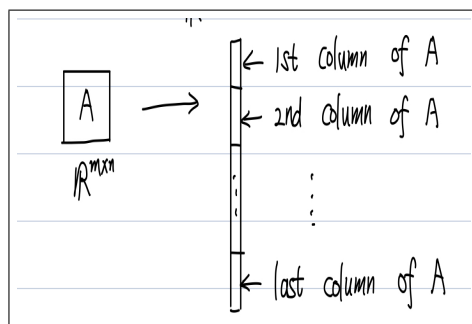
Geometric definition of different norms in \mathbb{R}^n



Note that $(\mathbb{R}^n, \|\cdot\|_1)$, $(\mathbb{R}^n, \|\cdot\|_2)$, $(\mathbb{R}^n, \|\cdot\|_\infty)$, ... are all different normed spaces. So for a given vector space, we can obtain various normed space by choosing different norms. Also, $\|x\|_p \leq \|x\|_q$ if $p \geq q$.

Example: $\mathbb{R}^{m \times n}$ is a vector space over \mathbb{R} .

1. $\mathbb{R}^{m \times n}$ can be viewed as \mathbb{R}^{mn} .



We can define vector p-norm for $\mathbb{R}^{m \times n}$.

- $p = 1$

$$\|A\|_{1,vec} = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|$$

- $p = 2$
 $\|A\|_{2,vec} = (\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2)^{\frac{1}{2}}$

This norm is widely known as the Frobenius norm denoted as $\|A\|_F$.

- $p = \infty$
 $\|A\|_{\infty,vec} = \max_{i=1,\dots,m} \max_{j=1,\dots,n} |a_{ij}|$

2. $\mathbb{R}^{m \times n}$ can be viewed as linear transformation from $\mathbb{R}^n \rightarrow \mathbb{R}^m$.
 We can define matrix p -norm for $\mathbb{R}^{m \times n}$.

$$\|A\|_p = \max_{x \neq 0, x \in \mathbb{R}^n} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \|Ax\|_p$$

- $p = 1$
 $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| = \text{maximum absolute column sum}$
- $p = \infty$
 $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| = \text{maximum absolute row sum}$
- $p = 2$
 $\|A\|_2 = \text{maximum singular value of } A$

3. We can also define other matrix norms.

- (a) We can use different norms in \mathbb{R}^n and \mathbb{R}^m .

$$\|A\|_{p \rightarrow q} = \max_{\|x\|_p=1} \|Ax\|_q$$

- (b) The nuclear norm $\|\cdot\|_*$

Example: $C[a, b]$ is a vector space over \mathbb{R} .

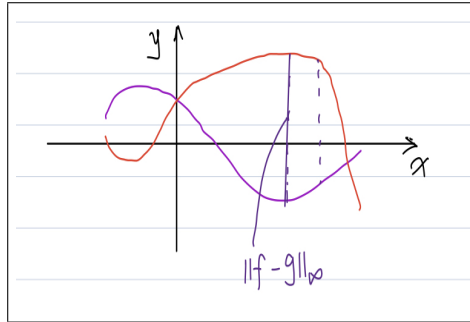
$\forall f \in C[a, b]$, define

$$\|f\|_\infty = \sup_{t \in [a, b]} |f(t)|$$

We can check that $\|\cdot\|_\infty$ is indeed a norm on $C[a, b]$.

The distance of two function $f, g \in C[a, b]$ is given by

$$\|f - g\|_\infty = \sup_{t \in [a, b]} |f(t) - g(t)|$$



Some other norms on $C[a, b]$.

1. $\|f\|_1 = \int_b^a |f(t)| dt$
2. $\|f\|_2 = (\int_b^a |f(t)|^2 dt)^{\frac{1}{2}}$
3. $\|f\|_p = (\int_b^a |f(t)|^p dt)^{\frac{1}{p}}$

Example: $L_\infty = \{a | a \text{ is a infinite sequence and } \exists c > 0 \text{ s.t. } |a_i| \leq c, \forall i\}$

1. $\forall a \in L_\infty$, define

$$\|a\|_\infty = \sup_i |a_i|$$

Remarks: You cannot replace sup here with max.

2. Define $\|a\|_p = (\sum_{i=1}^\infty |a_i|^p)^{\frac{1}{p}} \forall a \in L_\infty$ but this is not a norm on L_∞ .

$$\text{e.g. } a = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \vdots \\ \frac{1}{i} \\ \vdots \end{bmatrix} \in L_\infty, \text{ but } \|a\|_1 = \sum_{i=1}^\infty |a_i| = \sum_{i=1}^\infty \frac{1}{i} = \infty$$

So, $\|\cdot\|_1$ is not a norm on L_∞ .

Instead, we consider

$$L_p = \{a \in L_\infty | \|a\|_p < \infty\} \subset L_\infty$$

$$\|\cdot\|_p \text{ is a norm on } L_p.$$

$$\text{e.g. } a = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \vdots \\ \frac{1}{i} \\ \vdots \end{bmatrix} \in L_\infty$$

$$\|a\|_\infty = 1, \|a\|_2 = (\sum_{i=1}^\infty \frac{1}{i^2})^{\frac{1}{2}} = (\frac{\pi^2}{6})^{\frac{1}{2}} = \frac{\pi}{\sqrt{6}}, \|a\|_1 = \infty$$

So, $a \in L_\infty$, $a \in L_2$ but $a \notin L_1$. Indeed, $a \in L_p \forall p > 1$.

Limit and Convergence on Normed Vector Space

To define calculus, we first need to define convergent sequence.

Let \mathbb{V} be a normed vector space. Let $\{x^{(k)}\}_{k \in \mathbb{N}}$ be a sequence in \mathbb{V} , Let $x \in \mathbb{V}$. We say $\{x^{(k)}\}_{k \in \mathbb{N}}$ converges to x , denoted by $x^{(k)} \rightarrow x$, if

$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0$$
$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0 \iff x^{(k)} \rightarrow x$$

Example: Consider \mathbb{R}^n with $\|\cdot\|_2$,

$$\text{Let } x^{(k)} = \begin{bmatrix} \frac{1}{k} \\ \frac{2}{k} \\ \vdots \\ \frac{n}{k} \end{bmatrix} \in \mathbb{R}^n \text{ and } x = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^n$$

$$\|x^{(k)} - x\|_2 = \|x^{(k)}\|_2 = (\sum_{i=1}^n (\frac{i}{k})^2)^{\frac{1}{2}} = \frac{1}{k} (\sum_{i=1}^n i^2)^{\frac{1}{2}}$$
$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\|_2 = \lim_{k \rightarrow \infty} \frac{1}{k} (\sum_{i=1}^n i^2)^{\frac{1}{2}} = 0$$
$$x^{(k)} \rightarrow x$$

Unfortunately, the limit of a sequence may not always be in the same vector space as the original sequence. If this happen, we call this the normed vector space incomplete. Otherwise, it is a complete vector space also known as the Banach space.

Example of Banach space:

1. \mathbb{R}^n with any norm
2. $\mathbb{R}^{m \times n}$ with any norm
3. Tensor space $\mathbb{R}^{m \times n \times l}$ with any norm
4. $C[a, b]$ with $\|\cdot\|_\infty$
5. L_p with p-norm, for $p \geq 1$ and $p = \infty$.

Cauchy Sequence

Definition: $\{x^{(k)}\}$ is a Cauchy sequence, if for any $\epsilon > 0$, there exists K such that for any $k, l > K$, $\|x^{(k)} - x^{(l)}\| < \epsilon$.

Facts:

1. If $x^{(k)} \rightarrow x$ in $(\mathbb{V}, \|\cdot\|)$, then $\{x^{(k)}\}$ must also be a Cauchy sequence.

Proof.

$x^{(k)} \rightarrow x$ implies that $\forall \epsilon > 0, \exists k$, s.t. $k > K$ $\|x^{(k)} - x\| \leq \frac{\epsilon}{2}$. Therefore,
 $\|x^{(k)} - x^{(l)}\| \leq \|x^{(k)} - x\| + \|x^{(l)} - x\| \leq \epsilon, \forall k, l > K$

2. The reverse is **NOT** necessarily true.

Definition: A vector space $(\mathbb{V}, \|\cdot\|)$ is complete if the limit of all Cauchy sequences in \mathbb{V} is in \mathbb{V} .

Remarks: We can always complete an incomplete normed vector space by including all limits of its Cauchy sequence.

Finite Dimensional Vector Space

In most cases, we are dealing with finite dimensional vector space such as \mathbb{R}^n , $\mathbb{R}^{m \times n}$ and $\mathbb{R}^{m \times n \times l}$.

Properties related to Finite Dimensional Vector Space:

- Any finite dimensional vector space with any norm is complete. That is, any finite dimensional vector space is Banach space.
- For a finite dimensional vector space \mathbb{V} , all norms are equivalent.

Theorem: For any norms $\|\cdot\|_A$ and $\|\cdot\|_B$, $\exists c_1, c_2 > 0$ s.t.
 $c_1\|a\|_A \leq \|a\|_B \leq c_2\|a\|_A, \forall a \in \mathbb{V}$ (finite dimensional)

Example: Prove that $x^{(k)} \rightarrow x$ in $\|\cdot\|_A \iff x^{(k)} \rightarrow x$ in $\|\cdot\|_B$.
 Since $x^{(k)} \rightarrow x$ in $\|\cdot\|_A$,

$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\|_A = 0$$

Because of equivalence,

$$c_1\|x^{(k)} - x\|_A \leq \|x^{(k)} - x\|_B \leq c_2\|x^{(k)} - x\|_A$$

$$0 \leq \lim_{k \rightarrow \infty} \|x^{(k)} - x\|_B \leq c_2 \lim_{k \rightarrow \infty} \|x^{(k)} - x\|_A = 0$$

$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\|_B = 0 \text{ (by squeeze theorem)}$$

$$x^{(k)} \rightarrow x \text{ under } \|\cdot\|_B$$

Similarly for the \leftarrow direction.

Example: Consider \mathbb{R}^n and $\|\cdot\|_1$, $\|\cdot\|_2$ and $\|\cdot\|_\infty$.

- $\|\cdot\|_1$ and $\|\cdot\|_2$ are equivalent.

$$\|a\|_2 \leq \|a\|_1 \leq \sqrt{n}\|a\|_2, \forall a \in \mathbb{R}^n$$

- $\|\cdot\|_2$ and $\|\cdot\|_\infty$ are equivalent.

$$\|a\|_\infty \leq \|a\|_2 \leq \sqrt{n}\|a\|_\infty, \forall a \in \mathbb{R}^n$$

- $\|\cdot\|_1$ and $\|\cdot\|_\infty$ are equivalent.

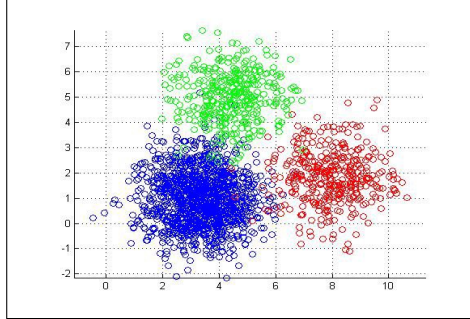
$$\|a\|_\infty \leq \|a\|_1 \leq n\|a\|_\infty, \forall a \in \mathbb{R}^n$$

Remarks: Though they are equivalent, the speed at which they converge are different. In other words, the convergence speed depends on norms.

Case Study: Clustering, k-means, k-medians

Clustering

Suppose we are given N vectors $x_1, x_2, \dots, x_N \in \mathbb{R}^n$, the goal of clustering is to group or partition the vectors into k groups or clusters, with the vectors in each group close to each other.



We use \mathbb{R}^n because it is simple, yet able to model a variety of data sets (e.g., signals, images, videos, attributes of things). Actually, the methods can be extended to any normed vector spaces (Banach space). Applications:

- Recommendation system
- Image clustering
- Text data clustering
- Many other applications.

Mathematical formulation:

- **Representation:**

Let $c_i \in \{1, 2, \dots, k\}$ be the group that x_i belongs to. $i = 1, 2, \dots, N$

Then, group G_j denoted by $G_j = \{i | c_i = j\}$. $j = 1, 2, \dots, k$.

We assign each group a representative vector, denoted by z_1, z_2, \dots, z_k .

The representative vectors are not necessarily one of the given vectors.

- **Evaluation:**

First of all, within one specific group G_j , all vectors should be close to the representative vector z_j . More precisely, let

$$J_j = \sum_{i \in G_j} \|x_i - z_j\|_2^2$$

Then, J_j should be small.

Secondly, consider all groups, since each J_j is small,

$$J = J_1 + J_2 + \dots + J_k$$

should be small.

Altogether, we solve the following

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} J \iff \min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k J_j \iff \min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k (\sum_{i \in G_j} \|x_i - z_j\|_2^2)$$

• **Optimization:**

We may use an alternating minimization to solve this minimization problem.

Step 1: Fix the representative z_1, \dots, z_k , find the best partitions G_1, \dots, G_k , i.e., solve

$$\min_{G_1, \dots, G_k} \sum_{j=1}^k (\sum_{i \in G_j} \|x_i - z_j\|_2^2) \quad \textcircled{1}$$

Step 2: Fix the groups G_1, \dots, G_k , find the best representatives z_1, \dots, z_k , i.e., solve

$$\min_{z_1, \dots, z_k} \sum_{j=1}^k (\sum_{i \in G_j} \|x_i - z_j\|_2^2) \quad \textcircled{2}$$

The two steps are repeated until convergence.

Let's find the solutions of the sub-problems $\textcircled{1}$ and $\textcircled{2}$ respectively.

For $\textcircled{1}$:

Finding the partition G_1, \dots, G_k is equivalent to finding c_1, \dots, c_N . So $\textcircled{1}$ becomes

$$\begin{aligned} \min_{c_1, \dots, c_N} (\|x_1 - z_{c_1}\|_2^2 + \dots + \|x_N - z_{c_N}\|_2^2) \\ \iff \min_{c_i \in \{1, 2, \dots, k\}} \|x_i - z_{c_i}\|_2^2 \\ \iff c_i = \operatorname{argmin}_{j \in \{1, 2, \dots, k\}} \|x_i - z_j\|_2^2 \end{aligned}$$

In other words, x_i is assigned to the group whose representative vector is the closest to x_i .

For $\textcircled{2}$:

It is rewritten as

$$\min_{z_1, \dots, z_k} (\sum_{i \in G_1} \|x_i - z_1\|_2^2 + \dots + \sum_{i \in G_k} \|x_i - z_k\|_2^2)$$

Obviously, it is equivalent to minimize each term independently, i.e., solve k independent problems.

$$\begin{aligned} \min_{z_j} (\sum_{i \in G_j} \|x_i - z_j\|_2^2), \quad j = 1, 2, \dots, k \\ \iff z_j = \frac{1}{|G_j|} \sum_{i \in G_j} x_i, \quad j = 1, 2, \dots, k \end{aligned}$$

where $|G_j|$ is the number of elements in G_j .

In other words, z_j is the mean of all vector in G_j . Note that in the above derivation, when we consider $n = 1$,

$$\begin{aligned} & \min_{z_j \in \mathbb{R}} \sum_{i \in G_j} (x_i - z_j)^2 \\ \iff & \sum_{i \in G_j} 2(x_i - z_j) = 0 \\ \iff & \sum_{i \in G_j} z_j = \sum_{i \in G_j} x_i \\ \iff & z_j = \frac{1}{|G_j|} \sum_{i \in G_j} x_i \end{aligned}$$

Altogether, we get the following clustering algorithm.

k-means Clustering

Initialization: Initialize z_1, z_2, \dots, z_k .

Step 1: Given z_1, z_2, \dots, z_k , compute

$$c_i = \underset{j \in \{1, 2, \dots, k\}}{\operatorname{argmin}} \|x_i - z_j\|_2^2, i = 1, 2, \dots, N$$

and define

$$G_j = \{i | c_i = j\}, j = 1, 2, \dots, k$$

Step 2: Given G_1, G_2, \dots, G_k , compute

$$z_j = \frac{1}{|G_j|} (\sum_{i \in G_j} x_i)$$

Go back to step 1.

In k-means, the Euclidean norm is used. We can replace it by 1-norm. We solve

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k \sum_{i \in G_j} \|x_i - z_j\|_1$$

k-medians Clustering

Initialization: Initialize z_1, z_2, \dots, z_k .

Step 1: Given z_1, z_2, \dots, z_k , compute

$$c_i = \underset{j \in \{1, 2, \dots, k\}}{\operatorname{argmin}} \|x_i - z_j\|_1, i = 1, 2, \dots, N$$

and define

$$G_j = \{i | c_i = j\}, j = 1, 2, \dots, k$$

Step 2: Given G_1, G_2, \dots, G_k , compute

$$z_j = \operatorname{median}\{x_i | i \in G_j\}$$

Go back to step 1.

2.3 Inner Product and Hilbert Space

Question: How do we describe the correlation/centerment between two vectors? Norms are not able to describe it as they are 'scaling sensitive'.

A good answer would be to use angle. A good candidate would be to use inner product since it is 'scaling insensitive'.

Inner Product

Definition: A function $\langle \cdot, \cdot \rangle : (\mathbb{V}, \mathbb{V}) \rightarrow \mathbb{R}$ on a vector space \mathbb{V} is called an inner product over \mathbb{R} , if:

1. $\forall x \in \mathbb{V}, \langle x, x \rangle \geq 0$ and $\langle x, x \rangle = 0 \iff x = 0$
2. $\langle \alpha x_1 + \beta x_2, y \rangle = \alpha \langle x_1, y \rangle + \beta \langle x_2, y \rangle, \forall \alpha, \beta \in \mathbb{R}, x_1, x_2, y \in \mathbb{V}$
3. $\langle x, y \rangle = \langle y, x \rangle, \forall x, y \in \mathbb{V}$

Remarks:

1. By 2 and 3, $\langle x, \alpha y_1 + \beta y_2 \rangle = \alpha \langle x, y_1 \rangle + \beta \langle x, y_2 \rangle, \forall \alpha, \beta \in \mathbb{R}, x_1, y_1, y_2 \in \mathbb{V}$. Therefore, $\langle \cdot, \cdot \rangle$ is a bi-linear function, i.e., it is linear with respect to one of the variable with the other fixed.
2. For inner product of vector spaces on \mathbb{C} , we only need to change 3 to $\langle x, y \rangle = \overline{\langle y, x \rangle}$, where $\langle \cdot, \cdot \rangle$ stands for complex conjugate.

Example: \mathbb{R}^n is a vector space. We can define an inner product as

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i = x^T y, \forall x, y \in \mathbb{R}^n.$$

Example: Another inner product in \mathbb{R}^n is as follows. We can define a "weighted" inner product as $\langle x, y \rangle_A = x^T A y$, where $A \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix.

Remarks: A is SPD $\iff A = A^T$ and $x^T A x > 0 \forall x \in \mathbb{R}^n$ and $x \neq 0$.

Example: $\mathbb{R}^{m \times n}$ is a vector space. We can define an inner product as

$$\langle A, B \rangle = \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ij}, \forall A, B \in \mathbb{R}^{m \times n}$$

Similarly, these are equal to $\text{trace}(A^T B)$, $\text{trace}(B^T A)$, $\text{trace}(AB^T)$ and $\text{trace}(BA^T)$, where $\text{trace}(A)$ is defined as the sum of the diagonal of matrix A .

Example: In L_2 , we can define an inner product as

$$\langle a, b \rangle = \sum_{i=1}^{\infty} a_i b_i, \forall a, b \in L_2$$

Example: In $C[a, b]$, we can define an inner product as

$$\langle f, g \rangle = \int_a^b f(t)g(t)dt, \forall f, g \in C[a, b]$$

Cauchy-Schwartz Inequality

If $\langle \cdot, \cdot \rangle$ is an inner product on \mathbb{V} , then, for any $x, y \in \mathbb{V}$,

$$|\langle x, y \rangle|^2 \leq \langle x, x \rangle \langle y, y \rangle$$

The equality holds true if and only if $x = \alpha y$ or $y = \alpha x$ for some $\alpha \in \mathbb{R}$

Proof.

Let $\lambda \in \mathbb{R}$ be an arbitrary number,

$$\begin{aligned} 0 &\leq \langle x + \lambda y, x + \lambda y \rangle \\ &= \langle x, x \rangle + \lambda \langle y, x \rangle + \lambda \langle x, y \rangle + \lambda^2 \langle y, y \rangle \\ &= \langle x, x \rangle + 2\lambda \langle x, y \rangle + \lambda^2 \langle y, y \rangle \end{aligned}$$

$$\text{Thus, } \lambda^2 \langle y, y \rangle + 2\lambda \langle x, y \rangle + \langle x, x \rangle \geq 0, \forall \lambda \in \mathbb{R}$$

The left is a quadratic function of λ and is always non-negative. There is at most one root of the quadratic function, hence, the determinant $b^2 - 4ac \leq 0$.

$$\text{So, } (2\langle x, y \rangle)^2 - 4\langle x, x \rangle \langle y, y \rangle \leq 0$$

$$\implies \langle x, y \rangle^2 \leq \langle x, x \rangle \langle y, y \rangle$$

Finally, when $\langle x, y \rangle^2 = \langle x, x \rangle \langle y, y \rangle$, there is a root, i.e., \exists a unique $\lambda \in \mathbb{R}$, $\lambda^2 \langle y, y \rangle + 2\lambda \langle x, y \rangle + \langle x, x \rangle = 0$.

$$\iff$$

$$\exists \text{ a unique } \lambda \in \mathbb{R}, \langle x + \lambda y, x + \lambda y \rangle = 0.$$

$$\iff$$

$$\text{a unique } \lambda \in \mathbb{R}, x + \lambda y = 0.$$

$$\iff$$

$$\exists \text{ a unique } \lambda \in \mathbb{R}, x = -\lambda y.$$

With the Cauchy-Schwartz inequality, we can show that

$$\|x\| = (\langle x, x \rangle)^{\frac{1}{2}} \text{ defines a norm.}$$

This is also called "norm induced by the inner product". This one above is for \mathbb{R}^n .

Proof.

$$\|x\| = (\langle x, x \rangle)^{\frac{1}{2}} \geq 0 \text{ and } \|x\| = (\langle x, x \rangle)^{\frac{1}{2}} = 0 \iff x = 0$$

$$\|\alpha x\| = (\langle \alpha x, \alpha x \rangle)^{\frac{1}{2}} = (\alpha^2 \langle x, x \rangle)^{\frac{1}{2}} = |\alpha| \|x\|$$

$$\begin{aligned} \|x + y\|^2 &= \langle x + y, x + y \rangle \\ &= \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle \\ &= \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle \\ &\leq \|x\|^2 + \|y\|^2 + 2\|x\|\|y\| \\ &= (\|x\| + \|y\|)^2 \\ \|x + y\| &\leq \|x\| + \|y\| \end{aligned}$$

Remarks: In the proof above, we have used an alternative version of the Cauchy-Schwartz inequality.

$$|\langle x, y \rangle| \leq \|x\|\|y\|$$

All kinds of induced norm

1. \mathbb{R}^n with inner product $\langle \cdot, \cdot \rangle : \langle x, y \rangle = x^T y$
The induced norm is
 $\|x\| = (\langle x, x \rangle)^{\frac{1}{2}} = (x^T x)^{\frac{1}{2}} = (\sum_{i=1}^n x_i^2)^{\frac{1}{2}} = \|x\|_2$
2. \mathbb{R}^n with weighted inner product $\langle \cdot, \cdot \rangle_A : \langle x, y \rangle_A = x^T A y$
The induced norm is
 $\|x\|_A = (x^T A x)^{\frac{1}{2}} = (\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j)$
3. The p-norm of \mathbb{R}^n
 $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}$
When $p = 2$, $\|\cdot\|_2$ is induced by $\langle \cdot, \cdot \rangle$. It is not induced by inner product for all p except for 2.
4. $\mathbb{R}^{m \times n}$ with inner product $\langle \cdot, \cdot \rangle : \langle A, B \rangle = \sum_{ij} a_{ij} b_{ij}$
The induced norm is
 $\|A\| = (\langle A, A \rangle)^{\frac{1}{2}} = (\sum_{ij} a_{ij}^2)^{\frac{1}{2}} = \|A\|_F = \|A\|_{vec, 2}$

5. Infinite sequence with inner product $\langle \cdot, \cdot \rangle : \langle a, b \rangle = \sum_{i=1}^{\infty} a_i b_i$
 $\|a\| = (\sum_{i=1}^{\infty} a_i^2)^{\frac{1}{2}} = \|a\|_2$
6. $C[a, b]$ with inner product $\langle \cdot, \cdot \rangle : \langle f, g \rangle = \int_a^b f(t)g(t)dt$
 $\|f\| = (\int_a^b (f(t))^2 dt)^{\frac{1}{2}} = \|f\|_2$

Angle in inner product spaces

By Cauchy-Schwartz inequality,

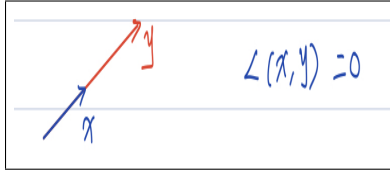
$$|\langle x, y \rangle| \leq \|x\| \|y\| \quad \forall x, y \in \mathbb{V}$$

Then,

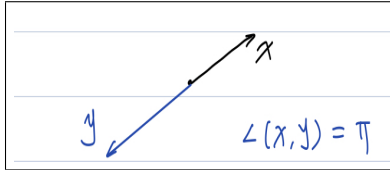
$$-\|x\| \|y\| \leq \langle x, y \rangle \leq \|x\| \|y\|$$

$$-1 \leq \frac{\langle x, y \rangle}{\|x\| \|y\|} \leq 1 \text{ if } x, y \neq 0$$

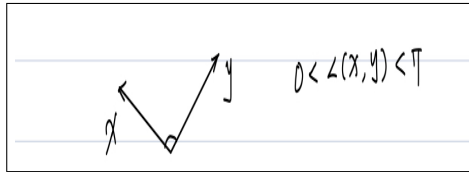
If $\frac{\langle x, y \rangle}{\|x\| \|y\|} = 1$, then $x = \alpha y$ with $\alpha > 0$. Otherwise, if $\alpha \leq 0$, then $\langle x, y \rangle = \alpha \langle y, y \rangle = \alpha \|y\|^2 \leq 0$. (*Contradiction*).



If $\frac{\langle x, y \rangle}{\|x\| \|y\|} = -1$, then $x = \alpha y$ with $\alpha < 0$.



If $-1 < \frac{\langle x, y \rangle}{\|x\| \|y\|} < 1$, then



Then we define

$$L(x, y) = \arccos \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

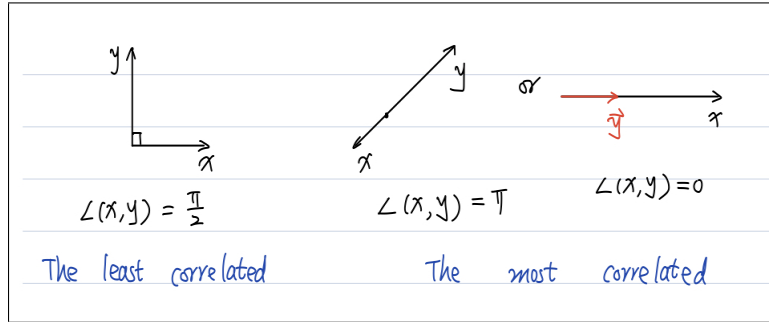
This definition is consistent with the observation above and the angles of vectors in \mathbb{R}^2 and \mathbb{R}^3 .

Orthogonality

Let \mathbb{V} be a vector space and $\langle \cdot, \cdot \rangle$ be the inner product.

- If $\frac{\langle x, y \rangle}{\|x\| \|y\|} = 1$ or -1 , then x and y are the most correlated.
- If $\frac{\langle x, y \rangle}{\|x\| \|y\|} = 0$, then x and y are the least correlated.

If $\langle x, y \rangle = 0$, then we say x and y are orthogonal.



Pythagorean theorem

Definition: Let x, y be two vectors in an inner product space \mathbb{V} .

Then $x \perp y \iff \|x + y\|^2 = \|x\|^2 + \|y\|^2$.

Proof.

$$\begin{aligned} \|x + y\|^2 &= \langle x + y, x + y \rangle \\ &= \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle \quad (1) \end{aligned}$$

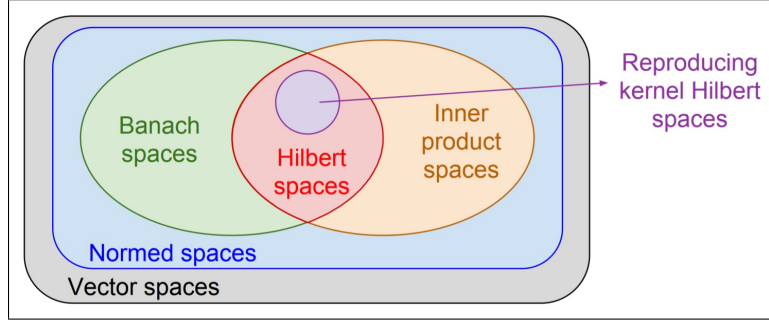
If $x \perp y$, then $\langle x, y \rangle = 0$.

$$\implies \|x + y\|^2 = \|x\|^2 + \|y\|^2$$

If $\|x + y\|^2 = \|x\|^2 + \|y\|^2$, together with (1), we have $\langle x, y \rangle = 0$.

Hilbert Space

Definition: A Hilbert space is a Banach space in which the norm is induced by an inner product.



Examples of Hilbert Space

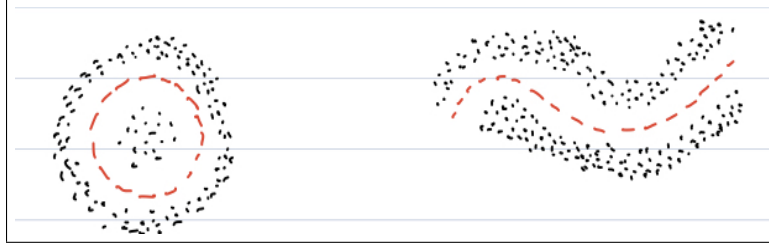
1. \mathbb{R}^n with $\langle \cdot, \cdot \rangle$ is a Hilbert space.
2. \mathbb{R}^n with $\langle \cdot, \cdot \rangle_A$ is a Hilbert space.
3. $\mathbb{R}^{m \times n}$ with $\langle \cdot, \cdot \rangle$ is a Hilbert space.
4. $L_2 = \{a \mid \|a\|_2 < \infty \text{ and } a \text{ is a infinite sequence}\}$ with $\langle \cdot, \cdot \rangle$ is a Hilbert space.
5. $C[a, b]$ with $\langle \cdot, \cdot \rangle$ is **NOT** a Hilbert space, because it is not a Banach space. In other words, the limit of a convergent sequence in $C[a, b]$ may not be in $C[a, b]$. To complete $C[a, b]$ under the norm $\|\cdot\| = (\langle \cdot, \cdot \rangle)^{\frac{1}{2}}$, we need to extend the Riemann integral to the so-called Lebesgue integral, and the resulting Hilbert space is $L^2(a, b)$.

In the following chapters, we will consider calculus on Hilbert/Banach spaces.

Case Study: Kernel trick, Kernel k-means

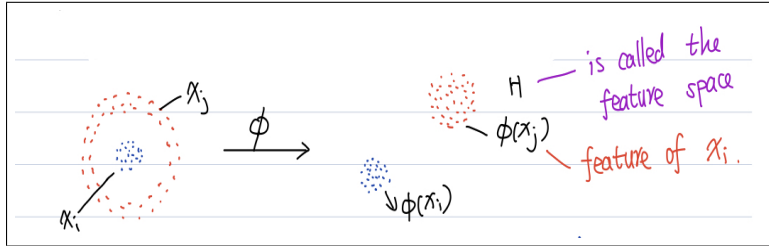
Recall that in k-means, we want to group $x_1, x_2, \dots, x_N \in \mathbb{R}^n$ into k groups. k-means work well only if the data are linearly separable. It will fail on "curved" data sets in \mathbb{R}^n .

k-means will fail for the following examples.



To modify k-means to these "curved" data sets in \mathbb{R}^n , we transform the "curved" data sets to "uncurved" data set in a Hilbert space \mathbb{H} .

Let $\phi : \mathbb{R}^n \rightarrow \mathbb{H}$,



Then we apply k-means to $\phi(x_1), \phi(x_2), \dots, \phi(x_N)$ in \mathbb{H} . Let z_1, z_2, \dots, z_k be the representative vectors in \mathbb{H} .

k-means Clustering

Step 0: Initialize z_1, z_2, \dots, z_k

Step 1: Given z_1, z_2, \dots, z_k , compute

$$c_i = \underset{j \in \{1, 2, \dots, k\}}{\operatorname{argmin}} \|\phi(x_i) - z_j\|^2, i = 1, 2, \dots, N$$

and define

$$G_j = \{i | c_i = j\}, j = 1, 2, \dots, k$$

Step 2: Given G_1, G_2, \dots, G_k , compute

$$z_j = \frac{1}{|G_j|} (\sum_{i \in G_j} \phi(x_i))$$

Go back to step 1.

However, finding the feature map ϕ is not easy, because ϕ depends on the shape of x_1, x_2, \dots, x_N , which generally is very complicated.

The good news is that:

There is no need to know ϕ explicitly in k-means algorithm.

Why?

- First of all, since we care only about the groups of x_1, \dots, x_N , we only need to know G_1, \dots, G_k . The representatives z_1, \dots, z_k are only intermediate. Therefore, we can eliminate z_1, \dots, z_k in the k-means algorithm.

Modified k-means Clustering

Step 0: Initialize G_1, \dots, G_k

Step 1: Given G_1, \dots, G_k , compute

$$c_i = \underset{j \in \{1, 2, \dots, k\}}{\operatorname{argmin}} \|\phi(x_i) - \frac{1}{|G_j|} \sum_{l \in G_j} \phi(x_l)\|^2, i = 1, 2, \dots, N$$

and define

$$G_j = \{i | c_i = j\}, j = 1, 2, \dots, k$$

Go back to step 1.

- Now, since we are talking about Hilbert space \mathbb{H} , we can expand the norm by

$$\begin{aligned} & \|\phi(x_i) - \frac{1}{|G_j|} \sum_{l \in G_j} \phi(x_l)\|^2 \\ &= \langle \phi(x_i) - \frac{1}{|G_j|} \sum_{l \in G_j} \phi(x_l), \phi(x_i) - \frac{1}{|G_j|} \sum_{l \in G_j} \phi(x_l) \rangle \\ &= \langle \phi(x_i), \phi(x_i) \rangle - \frac{2}{|G_j|} \sum_{l \in G_j} \langle \phi(x_i), \phi(x_l) \rangle + \\ & \quad \frac{1}{|G_j|^2} \sum_{l_1 \in G_j} \sum_{l_2 \in G_j} \langle \phi(x_{l_1}), \phi(x_{l_2}) \rangle \end{aligned}$$

All terms involved are in the form of

$$\langle \phi(x), \phi(y) \rangle : (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R}$$

Instead of defining ϕ explicitly, we define a function

$$\langle \phi(x), \phi(y) \rangle : (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R}$$

s.t. $k(x, y) = \langle \phi(x), \phi(y) \rangle$.

Therefore, an explicit expression of ϕ is **NOT** necessary. This process is also known as kernel trick.

Kernel function

$k(x, y)$ is called a kernel function. Which kernel function?

Necessary conditions:

$$1. k(x, y) = \langle \phi(x), \phi(y) \rangle = \langle \phi(y), \phi(x) \rangle = k(y, x)$$

We say k is a symmetric kernel is $k(x, y) = k(y, x), \forall x, y \in \mathbb{R}^n$.

$$2. \text{ Let } y_1, \dots, y_m \in \mathbb{R}^n. \text{ Then for any } C = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \in \mathbb{R}^n,$$

$$0 \leq \langle \sum_{i=1}^m c_i \phi(y_i), \sum_{j=1}^m c_j \phi(y_j) \rangle = \sum_{i=1}^m \sum_{j=1}^m c_i c_j \langle \phi(y_i), \phi(y_j) \rangle = \sum_{i=1}^m \sum_{j=1}^m c_i c_j k(y_i, y_j) = C^T k C \text{ where } k = [k(y_i, y_j)]_{i,j}$$

That is, $\forall C \in \mathbb{R}^m, C^T k C \geq 0$ and $k^T k$.

Definition: We say a kernel function $K : (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R}$ is symmetric positive semi-definite if:

$$1. k(x, y) = k(y, x), \forall x, y \in \mathbb{R}^n$$

2. For any m and $y_1, y_2, \dots, y_m \in \mathbb{R}^n$, the matrix

$$k = [k(y_i, y_j)]_{i,j}$$

is symmetric positive semi-definite.

Mercer's Theorem: If $K : (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R}$ is continuous and symmetric positive semi-definite, then there exists a Hilbert space \mathbb{H} and a mapping such that $k(x, y) = \langle \phi(x), \phi(y) \rangle$.

Some popular kernels:

- $k(x, y) = x^T y$ ($\phi(x) = x$ (No transform kernel))
- $k(x, y) = (x^T y + 1)^\alpha$ (α is an integer (Polynomial kernel))
- $k(x, y) = e^{-\frac{\|x-y\|_2^2}{\sigma^2}}$ ($\sigma > 0$ is a parameter (Gaussian kernel))

Kernel k-means Clustering

Step 0: Initialize G_1, \dots, G_k

Step 1: Given G_1, \dots, G_k , compute

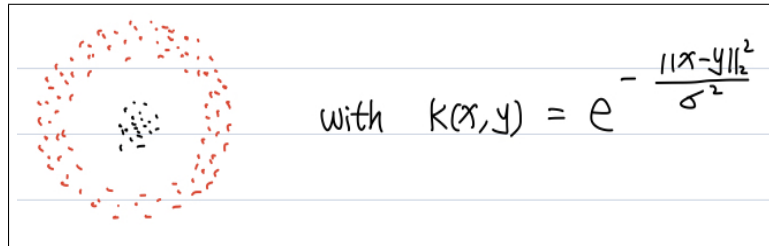
$$c_i = k(x_i, x_i) - \frac{2}{|G_j|} \sum_{l \in G_j} k(x_i, x_l) + \frac{1}{|G_j|^2} \sum_{l_1 \in G_j} \sum_{l_2 \in G_j} k(x_{l_1}, x_{l_2})$$

and define

$$G_j = \{i | c_i = j\}, j = 1, 2, \dots, k$$

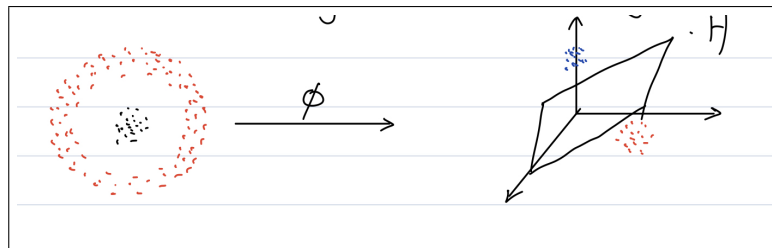
Go back to step 1.

Why kernel k-means work?



Suppose we use Gaussian kernel,

- $k(x_i, x_i) = e^{-\frac{\|x_i - x_i\|_2^2}{\sigma^2}} = e^{-0} = 1, \forall i$
so all $\phi(x_1), \dots, \phi(x_N)$ are on unit sphere in H .
- $k(x_i, x_j) \begin{cases} \approx 1 \text{ if } x_i \approx x_j \\ \approx 0 \text{ if } \|x_i - x_j\|_2 \text{ is large} \end{cases}$ Since $\langle \phi(x_i), \phi(x_j) \rangle = k(x, y)$,
 - If $x_i \approx x_j$, then
 $\|\phi(x_i) - \phi(x_j)\|^2 = \|\phi(x_i)\|^2 - 2\langle \phi(x_i), \phi(x_j) \rangle + \|\phi(x_j)\|^2 \approx 0$
 $\implies \phi(x_i) = \phi(x_j)$.
 - If $\|x_i - x_j\|_2$ is large, then $\phi(x_i) \perp \phi(x_j)$.

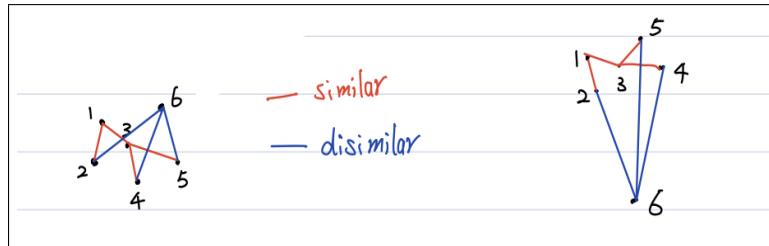


Thus, kernel k-means work for "curved" data sets which k-means fail.

Case Study: Metric Learning

Given a set of data $x_1, x_2, \dots, x_n \in \mathbb{R}^n$ and $S : (x_i, x_j) \in S$ if x_i and x_j are similar and $D : (x_i, x_j) \in D$ if x_i and x_j are dissimilar.

Our goal is to find a "new" metric such that for similar pair, it is close and for dissimilar pair, it is far away. In other words, in metric learning, given a set of data in two groups, we need to find the metric that would differentiate between the two groups.



Representation:

Norm induced by weighted inner product

Given $A \in \mathbb{R}^{n \times n}$ is SPD,

$$\langle x, y \rangle = x^T A y$$

$$\text{and } \|x\|_A = (x^T A x)^{\frac{1}{2}}$$

Then finding a metric is the same as finding an SPD matrix A .

Remarks:

1. The set of all SPD is **NOT** closed.
2. The closure of the set of all SPD matrices is the set of all SPSD matrices.
3. If A is SPSD, then $\|x\|_A$ is not a norm because $\|x\|_A^2 = 0 \iff x^T A x = 0$ cannot implies $x = 0$.
4. $\|\cdot\|_A$ is still a semi-norm:
 - $\|x\|_A \geq 0$
 - $\|\alpha x\|_A = |\alpha| \|x\|_A$
 - $\|x + y\|_A \leq \|x\|_A + \|y\|_A$

Evaluation:

Which A is the best?

1. For $(x_i, x_j) \in S$, $dist(x_i, x_j)$ should be small

$$\sum_{(x_i, x_j) \in S} \|x_i - x_j\|_A^2$$

2. For $(x_i, x_j) \in D$, $dist(x_i, x_j)$ should not be small

Altogether:

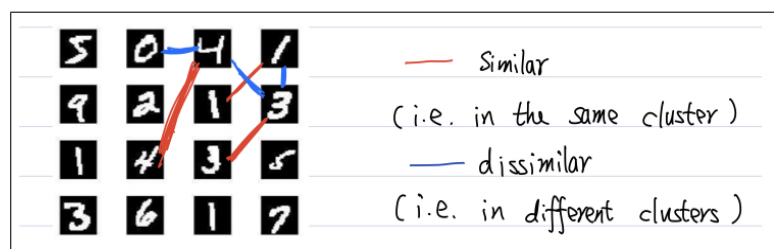
$$\begin{aligned} \min_{A \in \mathbb{R}^{n \times n}} \quad & \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_A^2 \\ \text{s.t.} \quad & \sum_{(x_i, x_j) \in D} \|x_i - x_j\|_A^2 \geq 1 \end{aligned}$$

Optimization:

It is too complicated.

One popular application: "Supervised" clustering

- Given a data set with partial clustering information,



- Apply metric learning (i.e. find a distance)
- Cluster the points under the newly learned distance metric.

Linear and Differentiable Functions

3.1 Linear Function

Definition: Let \mathbb{V} be a vector space and $f : \mathbb{V} \rightarrow \mathbb{R}$ be a function. f is a linear function if

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y), \forall \alpha, \beta \in \mathbb{R} \text{ and } x, y \in \mathbb{V}$$

Example: The mean of a vector in \mathbb{R}^n .

$$\forall x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n, f(x) = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n} \text{ is a linear function because}$$

$$f(\alpha x + \beta y) = \frac{\sum_{i=1}^n (\alpha x_i + \beta y_i)}{n} = \alpha \frac{\sum_{i=1}^n x_i}{n} + \beta \frac{\sum_{i=1}^n y_i}{n} = \alpha f(x) + \beta f(y)$$

Example: The maximum entry of a vector in \mathbb{R}^n .

$$\forall x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n, f(x) = \max_{i=1, \dots, n} x_i \text{ is not a linear function.}$$

One counter example:

$$x = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \alpha = 1, \beta = 1$$
$$f(\alpha x + \beta y) = f\left(\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}\right) = 1 \text{ but } \alpha f(x) = 1 \text{ and } \beta f(y) = 1$$

Hence, $f(\alpha x + \beta y) \neq \alpha f(x) + \beta f(y)$.

Example: $f : \mathbb{R}^n \rightarrow \mathbb{R}$ defined by $f(x) = \langle a, x \rangle$, where $a \in \mathbb{R}^n$ is a fixed vector in \mathbb{R}^n is linear.

Example: $F : C[-1, 1] \rightarrow \mathbb{R}$ defined by $F(f) = f(0)$ is linear because

$$F(\alpha f + \beta g) = (\alpha f + \beta g)(0) = \alpha f(0) + \beta g(0) = \alpha F(f) + \beta F(g)$$

Example: $F : C[a, b] \rightarrow \mathbb{R}$ defined by $F(f) = \int_a^b f(t)dt$ is linear because

$$\begin{aligned} F(\alpha f + \beta g) &= \int_a^b (\alpha f + \beta g)(t)dt \\ &= \int_a^b (\alpha f(t) + \beta g(t))dt \\ &= \alpha \int_a^b f(t)dt + \beta \int_a^b g(t)dt \\ &= \alpha F(f) + \beta F(g) \end{aligned}$$

Example: Let \mathbb{V} be an inner product space with inner product $\langle \cdot, \cdot \rangle$. Let $a \in \mathbb{V}$ and $f : \mathbb{V} \rightarrow \mathbb{R}$ defined by $f(x) = \langle a, x \rangle$ is linear.

Example: A norm function on \mathbb{V} is **NOT** linear.

Proof: Let $\| \cdot \| : \mathbb{V} \rightarrow \mathbb{R}$. Then $\| -x \| = \|x\|$ by norm property. If $\| \cdot \|$ is linear, then

$$\| -x \| = \| -x + 0 \cdot x \| = -1\|x\| + 0\|x\| = -\|x\| \text{ (Contradiction)}$$

Properties of Linear Function:

1. *Homogeneity*:

$$f(\alpha x) = \alpha f(x), \forall \alpha \in \mathbb{R}, x \in \mathbb{V}$$

$$\text{because } f(\alpha x) = f(\alpha x + 0 \cdot y) = \alpha f(x) + 0 \cdot f(y) = \alpha f(x)$$

Choosing $\alpha = 0$, then we obtain $f(0) = 0$.

2. *Additivity*:

$$f(x + y) = f(x) + f(y)$$

$$f(\alpha_1 x_1 + \alpha_2 x_2 + \cdots + \alpha_k x_k)$$

$$= \alpha_1 f(x_1) + f(\alpha_2 x_2 + \cdots + \alpha_k x_k)$$

$$= \alpha_1 f(x_1) + \alpha_2 f(x_2) + f(\alpha_3 x_3 + \cdots + \alpha_k x_k)$$

$$= \cdots$$

$$= \alpha_1 f(x_1) + \alpha_2 f(x_2) + \cdots + \alpha_k f(x_k)$$

Linear Function on Hilbert Space

For simplicity, let's consider a linear function on \mathbb{R}^n equipped with the standard inner product $\langle x, y \rangle = x^T y$ and the induced norm $\|x\|_2 = (\langle x, x \rangle)^{\frac{1}{2}}$.

- From one of the examples above,

For any give $a \in \mathbb{R}^n$, the function $f(x) = \langle a, x \rangle$ is linear.

- The reverse is true, i.e.,

Any linear function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ must be in the form of $f(x) = \langle a, x \rangle$ for some $a \in \mathbb{R}^n$.

We are assuming that $\mathbb{H} = \mathbb{R}^n$ for simplicity, but this theorem actually holds for any forms of Hilbert Space \mathbb{H} .

Theorem: For any linear function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, there exists a unique $a \in \mathbb{R}^n$ s.t. $f(x) = \langle a, x \rangle$, $\forall x \in \mathbb{R}^n$.

Proof: Let e_1, e_2, \dots, e_n be the natural basis of \mathbb{R}^n where e_i is a vector where the i -th entry is 1 and 0 elsewhere.

$$\forall x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n, x = x_1 e_1 + x_2 e_2 + \dots + x_n e_n$$

$$\text{So } f(x) = f(x_1 e_1 + x_2 e_2 + \dots + x_n e_n)$$

$$= x_1 f(e_1) + x_2 f(e_2) + \dots + x_n f(e_n)$$

$$= \left\langle \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} f(e_1) \\ f(e_2) \\ \vdots \\ f(e_n) \end{bmatrix} \right\rangle$$

$$= \langle a, x \rangle \text{ where } a = \begin{bmatrix} f(e_1) \\ f(e_2) \\ \vdots \\ f(e_n) \end{bmatrix}$$

Now we prove the uniqueness of this theorem.

Suppose a is **NOT** unique, $\exists a, b \in \mathbb{R}^n$ s.t.

$$f(x) = \langle a, x \rangle = \langle b, x \rangle, \forall x \in \mathbb{R}^n$$

Then choose $x = e_i, i = 1, \dots, n$

$$f(e_i) = \langle a, e_i \rangle = \langle b, e_i \rangle \implies a_i = b_i, i = 1, \dots, n$$

$$\implies a = b$$

(Contradiction)

Therefore, a must be unique.

Riesz Representation Theorem

Extending previous theorem to the entirety of Hilbert space \mathbb{H} .

Theorem:

Let \mathbb{H} be a Hilbert space. Let $f : \mathbb{H} \rightarrow \mathbb{R}$. Then f is linear and bounded if and only if $f(x) = \langle a, x \rangle$ for some unique $a \in \mathbb{H}$.

Example: We know that mean of a vector on \mathbb{R}^n is linear.

$$f(x) = \frac{x_1 + x_2 + \dots + x_n}{n} = \left\langle \frac{1}{n} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, x \right\rangle$$

Example: Let \mathbb{H} be a Hilbert space and $\|\cdot\|$ is **NOT** linear. So there is no such $a \in \mathbb{H}$ s.t. $\|x\| = \langle a, x \rangle, \forall x \in \mathbb{H}$.

Example: $\mathbb{R}^{n \times n}$ with inner product

$$\langle A, B \rangle = \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ij}, \forall A, B \in \mathbb{R}^{n \times n}$$

Define $trace(A) = \sum_{i=1}^n a_{ii}$, $\forall A \in \mathbb{R}^{n \times n}$ is linear. We have

$$trace(A) = \langle A, I \rangle$$

Remarks:

1. In finite dimensional Hilbert space, linear \iff linear and bounded.
2. In infinite dimensional Hilbert space, there exists linear but unbounded function.

Example: $L^2(-1, 1)$ - the completion of $C[-1, 1]$ under the inner product $\langle f, g \rangle = \int_{-1}^1 f(t)g(t)dt$ and $\|f\|_2 = (\langle f, f \rangle)^{\frac{1}{2}} = (\int_{-1}^1 |f(t)|^2 dt)^{\frac{1}{2}}$. Consider $F(f) = f(0)$, $\forall f \in L^2(-1, 1)$
But $F(f)$ is unbounded since

$$\exists f \in L^2(-1, 1) \text{ s.t. } F(f) = \infty$$

$$\text{e.g. } f(t) = \begin{cases} 1 & t \neq 0 \text{ and } t \in (-1, 1) \\ \infty & t = 0 \end{cases}$$

There exists no inner product representation for $F(f) = f(0)$.

Example: $L^2(-1, 1)$
Consider $G : L^2(-1, 1) \rightarrow \mathbb{R}$

$$G(f) = \int_{-1}^1 f(t)dt$$

G is linear.

G is bounded because for any $f \in L^2(-1, 1)$,

$$G(f) = \int_{-1}^1 f(t)dt = \int_{-1}^1 f(t) \cdot 1dt = \langle f, 1 \rangle \leq \|f\|_2 (\int_{-1}^1 1^2 dt)^{\frac{1}{2}} \leq 2\|f\|_2$$

Riesz $\implies g \in L^2(-1, 1)$ s.t. $G(f) = \langle f, g \rangle$. Indeed, $g(t) = 1$, $\forall t \in (-1, 1)$.

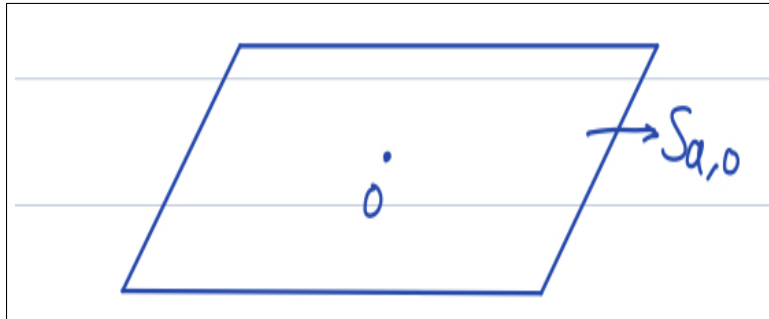
Hyperplane

Let \mathbb{H} be a Hilbert space and $a \in \mathbb{H}$.

Consider $S_{a,0} = \{x \in H | \langle a, x \rangle = 0\} \subset \mathbb{H}$,

Then $\forall \alpha, \beta \in \mathbb{R}$ and $\forall x, y \in S_{a,0}$

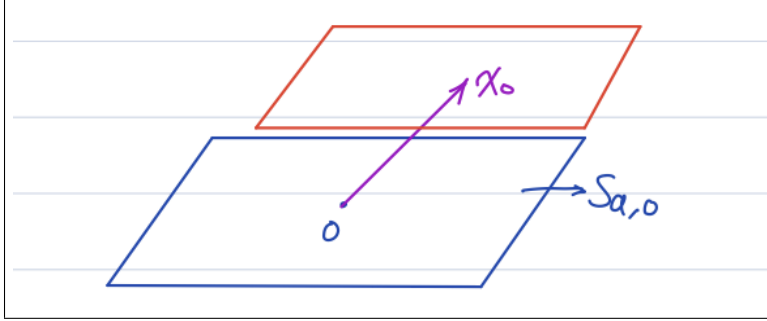
$\langle a, \alpha x + \beta y \rangle = \alpha \langle a, x \rangle + \beta \langle a, y \rangle = 0$. That is, $\alpha x + \beta y \in S_{a,0} \implies S_{a,0}$ is a linear space (subspace of H).



$S_{a,b} = \{x \in \mathbb{H} | \langle a, x \rangle = b\} \subset \mathbb{H}$
Let $x_0 \in S_{a,b}$, then $\langle a, x_0 \rangle = b$.

1. $\forall x \in S_{a,b}$
 $\langle a, x - x_0 \rangle = \langle a, x \rangle - \langle a, x_0 \rangle = b - b = 0$
 $\implies x - x_0 \in S_{a,0} \implies x \in x_0 + S_{a,0} \implies S_{a,b} \subset x_0 + S_{a,0}$
2. $\forall x \in S_{a,0}$
 $\langle a, x + x_0 \rangle = \langle a, x \rangle + \langle a, x_0 \rangle = 0 + b = b$
 $\implies x + x_0 \in S_{a,b} \implies S_{a,0} + x_0 \subset S_{a,b}$

(1) and (2) $\implies S_{a,b} = S_{a,0} + x_0$
 $S_{a,b}$ is a shift of a subspace.



Thus, $S_{a,b}$ is a plane on \mathbb{H} . So we call $S_{a,b}$ a hyperplane. Also, its co-dimension is 1 because it is defined by one linear equation.

Projection Onto Hyperplane

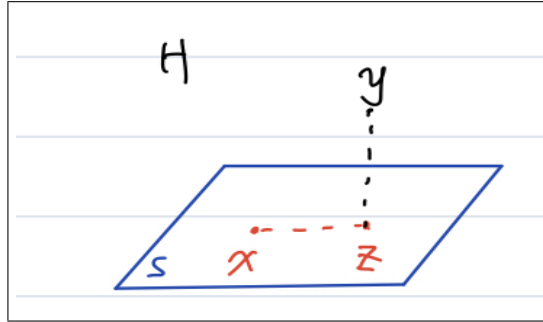
Consider a hyperplane in \mathbb{H}

$$S = \{x \in \mathbb{H} | \langle a, x \rangle = b\}$$

Given $y \in \mathbb{H}$, the vector on S that is closest to y is called the projection of y onto S , denoted by $P_S y$.

$$P_S y = \operatorname{argmin}_{x \in S} \|y - x\|$$

Our goal here is to find the explicit form of $P_S y$ in terms of a , b and y .



Theorem: z is a solution of $\operatorname{argmin}_{x \in S} \|y - x\| \iff z \in S$ and $\langle z - y, x - z \rangle = 0$, $\forall x \in S$.

Remarks: Since $\forall x \in S, x - z \in S - z$ and $z \in S \implies x - z \in S_{a,0}$, so $\langle z - y, x - z \rangle = 0$ implies $z - y \perp S_{a,0}$.

Proof: (\implies) Assume z is a solution of $\operatorname{argmin}_{x \in S} \|y - x\|$, then $z \in S$. $\forall x \in S$ and $\forall t \in \mathbb{R}$,

$$\langle a, z + t(x - z) \rangle = \langle a, z \rangle + t(\langle a, x \rangle - \langle a, z \rangle) = b + t(b - b) = b$$

Hence, $z + t(x - z) \in S$.
Since z is a minimizer,

$$\begin{aligned} \|z - y\|^2 &\leq \|z + t(x - z) - y\|^2 = \|z - y + t(x - z)\|^2 \\ &= \|z - y\|^2 + 2t\langle z - y, x - z \rangle + t^2\|x - z\|^2 \\ &\implies 2t\langle z - y, x - z \rangle \geq -t^2\|x - z\|^2 \end{aligned}$$

- If $t > 0$, then

$$\langle z - y, x - z \rangle \geq -\frac{t}{2}\|x - z\|^2$$

Let $t \rightarrow 0_+$, then

$$\langle z - y, x - z \rangle \geq \lim_{t \rightarrow 0_+} (-\frac{t}{2}\|x - z\|^2) = 0$$

- If $t < 0$, then

$$\langle z - y, x - z \rangle \leq -\frac{t}{2}\|x - z\|^2$$

Let $t \rightarrow 0_-$, then

$$\langle z - y, x - z \rangle \leq \lim_{t \rightarrow 0_-} (-\frac{t}{2}\|x - z\|^2) = 0$$

$$\implies \langle z - y, x - z \rangle = 0$$

(\Leftarrow) Assume $z \in S$ and $\langle z - y, x - z \rangle = 0, \forall x \in S$,

$$\begin{aligned} \|x - y\|^2 &= \|(x - z) + (z - y)\|^2 \\ &= \|x - z\|^2 + 2\langle x - z, z - y \rangle + \|z - y\|^2 \\ &= \|x - z\|^2 + \|z - y\|^2 \geq \|z - y\|^2 \\ &\implies z = \operatorname{argmin}_{x \in S} \|x - y\| \end{aligned}$$

Theorem: Let \mathbb{H} be a Hilbert space and $a \in \mathbb{H}$. Let $b \in \mathbb{R}$ and $S = \{x \in \mathbb{H} | \langle a, x \rangle = b\}$. Given $y \in \mathbb{H}$, the solution of

$$\operatorname{argmin}_{x \in S} \|y - x\|$$

exists and is unique, which is given by

$$y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$$

Proof: $z = y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$

Then

$$\begin{aligned} 1. \quad \langle a, z \rangle &= \langle a, y \rangle - \langle a, \frac{\langle a, y \rangle - b}{\|a\|^2} a \rangle \\ &= \langle a, y \rangle - \frac{\langle a, y \rangle - b}{\|a\|^2} \langle a, a \rangle \\ &= \langle a, y \rangle - (\langle a, y \rangle - b) \\ &= b \implies z \in S \\ 2. \quad \text{For any } x \in S, \\ \langle z - y, x - z \rangle &= \langle (-\frac{\langle a, y \rangle - b}{\|a\|^2}) a, x - z \rangle \\ &= -\frac{\langle a, y \rangle - b}{\|a\|^2} (\langle a, x \rangle - \langle a, z \rangle) \\ &= -\frac{\langle a, y \rangle - b}{\|a\|^2} (b - b) \\ &= 0 \end{aligned}$$

Hence, z is a solution of $\operatorname{argmin}_{x \in S} \|y - x\|$. It remains to check the uniqueness.

Suppose it has two solutions z_1 and z_2 . Then $z_1, z_2 \in S$.

$$\begin{aligned} z_1 \text{ is a solution} &\implies \langle z_1 - y, z_2 - z_1 \rangle = 0 \\ z_2 \text{ is a solution} &\implies \langle z_2 - y, z_1 - z_2 \rangle = 0 \implies \langle y - z_2, z_2 - z_1 \rangle = 0 \end{aligned}$$

Adding the two identities, we have

$$\begin{aligned} \langle z_1 - z_2, z_2 - z_1 \rangle &= 0 \\ \iff -\|z_1 - z_2\|^2 &= 0 \\ \iff z_1 &= z_2 \text{ (Contradiction)} \end{aligned}$$

In summary,

Let \mathbb{H} be a Hilbert space and

$$S = \{x \in \mathbb{H} | \langle a, x \rangle = b\}$$

Let $y \in \mathbb{H}$. Then the projection of y onto S is

$$P_S y = \operatorname{argmin}_{x \in S} \|y - x\|$$

and can be given by

$$P_S y = y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$$

Affine Function

Definition: A linear function plus a constant is an affine function. i.e. f is affine if $f(x) = g(x) + b$ where $g : \mathbb{H} \rightarrow \mathbb{R}$ is linear and $b \in \mathbb{R}$.

Properties:

1. If $f : \mathbb{H} \rightarrow \mathbb{R}$ is affine, then for any $\alpha, \beta \in \mathbb{R}$ and $\alpha + \beta = 1$, we have

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$$

To see this,

$$\begin{aligned} f(\alpha x + \beta y) &= g(\alpha x + \beta y) + (\alpha + \beta)b \\ &= \alpha g(x) + \beta g(y) + (\alpha + \beta)b \\ &= \alpha(g(x) + b) + \beta(g(y) + b) \\ &= \alpha f(x) + \beta f(y) \end{aligned}$$

2. If \mathbb{H} is a Hilbert space and f is bounded, then f is affine if and only if

$$f(x) = \langle a, x \rangle + b \text{ for some } a \in \mathbb{H} \text{ and } b \in \mathbb{R}$$