Ch. 2. Normal Spaces / Inner product spaces

§ 2.1. Vector spaces · Definition: A vector space over IR (the real domain) is a set 1/ together with two functions: vector addition: $f: V \times V \rightarrow V$ (i.e., $\chi + y$, where χ, y) scalar multiplication: ·: RXV -> V (i.e., XX, where XER, XEV) that satifying the following. (D) Associativity of addition: X + (Y+Z) = (X+Y) + Z $\forall X, Y, Z \in V$ (2) Commutativity of addition: X+Y = Y+X $\forall X, Y \in V$. 3 Zero vector: I an element, denoted by 0, in V, s.t. $\chi_{+0} = \vartheta + \chi \simeq \chi \qquad \forall \ \chi \in V.$ 4 Negative vector: YXEV, I an element, denoted by -X, in V, S.t. $\chi + (-\chi) = (-\chi) + \chi = 0$ (5) $\forall x \in V, \quad |x = x.$ $\forall x \in V \text{ and } \forall, \beta \in \mathbb{R}, \quad \alpha(\beta x) = (\alpha \beta) x$. (6) $\forall x \in V \text{ and } x, \beta \in \mathbb{R}$, $(\alpha + \beta) x = \alpha x + \beta x$ 67) YX, yEV and XER, X(X+y) = XX+XY @ 8 Remark: • We can define vector space over ((the complex domain) similarly. · We will assume vector space over IR for default. Vector space over (is used very rarely.

Example 1: R is a vector space, with "+" the standard addition of real numbers and "•" the standard multiplication of real numbers. Example 2: R" is a vector space, with "+" and "." defined by:

addition: $\forall \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in I\mathbb{R}^n, \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} X_1 + y_1 \\ X_2 + y_2 \\ \vdots \\ X_n + y_n \end{bmatrix}$
scalar multiplication: $\forall \ x \in \mathbb{R} \text{ and } \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n, x \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \alpha \times x_1 \\ \alpha \times x_2 \\ \vdots \\ \alpha \times x_n \end{bmatrix}.$
Many input data can be modeled by vectors in IR ⁿ .
 Digital sound signals of length n.
• Time series of length n.
• n different attributes/features of a single thing or object.
Example 3: All real mxn matrices is a vector space, with standard
matrix addition and standard scalar multiplication.
· This vector space is the same as IR ^{mn} .
. An mxn matrix can be used to represent a black-white
digital image of mxn pixels.
Example 4: All 3-arrays of size mxnxl is a vector space,
if "t" and "." is defined by
+: $\forall x = [x_{ijk}]_{i=1}^{n} \underset{j=1}{\overset{m}{k}} x_{i=1} and y = [y_{ijk}]_{i=1}^{n} \underset{j=1}{\overset{m}{k}} x_{i=1}$
$\chi + y = [\chi_{ijk} + y_{ijk}]_{i=1}^{n} + \frac{m}{1-1} + \frac{1}{k-1}$
• $\forall x = [x_{ijk}]_{i=1}^n \stackrel{m}{=} k and x \in \mathbb{R}$
$X = [X \times i]_{i=1}^{n} = K_{i}$
• This vector space is the same as IR mal
· An nxmx3 3-among can be used to model a color digital
image, where Xijk means (iij) the pixel in channel K, and
Channel 1, 2, 3 means Red, Green, Blue channels of the image.
• An nxmxl 3-array can be used to model a video,
where Xijk means the (i,j)-th pixel at K-th frame.

Example 5: Consider the set of all strings.
Define the addition by, e.g.,
'I ' + 'an ' = 'I an '
and some scalar multiplication.
Then it doesn't form a vector space.
• Therefore, we cannot use vector space to model text data in
this naive way.
Example 6: The function space C[a,b] = { f f is continous on [a,b] }
is a vector space if we define "+" and "." by:
$+$: \forall f, g \in C[a, b], (f+g)(t) = f(t)+g(t), \forall t $\in [a, b]$.
• : $\forall f \in C[a,b], \forall \in \mathbb{R}, (\forall f)(t) = \forall f(t))$
·C[a,b] is referred to as a function space, since any vector
in the vector space is a function.
· C[a, b] could be the hypothesis space of a learner with
one input and one output, i.e.,
$\chi_i \rightarrow Y_i$, with $\chi_i \in [a, b]$ and $Y_i \in R$.
leave a $f\in CIA,b]$ s.t. $f(X_i) \approx Y_i$ for all i.

§ 1.2. Normed spaces and Banach Spaces
In order to do calculus on vector spaces, we need to
define 'distance, close ness between vectors.
Let V be a vector space. Let
$$X, Y \in V$$
. Then,
distance $(X, Y) = distance (X - Y, Y - Y) = distance (X + Y, 0)$
(distance should be shift invariant) length of X-Y.
Therefore, to define a distance, we only need to define
a length for each vector in V.
Let $X \in V$. Let $||X||$ be its length, called norm, which
should satisfy
(D a length should be non negative, re.
 $||X|| \ge 0$ $\forall X \in V$.
Moreover, only the zero vector has a zero length, re,
 $||X|| \ge 0$ $\forall X \in V$.
Moreover, only the zero vector should be
the multiple of the length of the vector. i.e.,
 $\forall X \in \mathbb{R}$. $||X|| = |X|| ||X||$
(3) Triangular inequality: the length of the direct path is
the smallest
 $||X+Y|| \le ||X|| + ||Y||$
 $||X+Y|| \le ||X|| + ||Y||$

As p-norm
$$(p \ge 1)$$

 $\|X\|_p = \left(\sum_{i=1}^{n} |X_i|^p\right)^{h_p}$
• Comparison of Unit balls.
 $(x_i|x_i|_p=1)$
 $(x_i|_p=1)$
 $($

Some other norms of C[a,b] can be

$$\|f\|_{1} = \int_{a}^{b} |f(x)| dx$$

$$\|f\|_{2} = \left(\int_{a}^{b} |f(x)|^{2} dx\right)^{\frac{1}{2}}$$

$$\|f\|_{p} = \left(\int_{a}^{b} |f(x)|^{p} dx\right)^{\frac{1}{p}}$$

To define calculus, we need first define convergent sequence.
Let V be a normed vector space.
Let
$$\{X_k\}_{k\in\mathbb{N}}$$
 be a sequence in V, (i.e., $X_k \in V$ $\forall k = 1, 2, 3, \cdots$).
Let $\{X_k\}_{k\in\mathbb{N}}$ be a sequence in V, (i.e., $X_k \in V$ $\forall k = 1, 2, 3, \cdots$).
Let $X \in V$. We say $\{X_k\}_{k\in\mathbb{N}}$ converges to X, denoted by $X_k \to X$, if
 $\lim_{k \to \infty} ||X_k - X|| = 0$
Example 1: Consider $|\mathbb{R}^n$ with $||\cdot||_2$.
Let $X_k = \binom{v_k}{v_k} \operatorname{GIR}^n \forall k$ and $X = \binom{0}{i} = 0$
Then, $||X_k - X||_2 = ||\binom{v_k}{v_k} \operatorname{GIR}^n \forall k$ and $X = \binom{0}{i} = 0$
Then, $||X_k - X||_2 = ||\binom{v_k}{v_k} \operatorname{GIR}^n = F(n) \cdot \binom{im}{k} = \frac{F(n)}{n}$.
 $\lim_{k \to \infty} ||X_k - X||_2 = \lim_{k \to \infty} \frac{F(n)}{k} = F(n) \cdot \binom{im}{k} = 0$.
Therefore, $X_k \to X$ as $k \to \infty$.
Example 2: Consider C[0, 1] with $||\cdot||_\infty$
Let $f_k(t) = \operatorname{Sin}(2\pi kt)/k^2$
Then, $||f_k - 0||_\infty = \sup_{k\to\infty} \int_{k}^{\sin} (2\pi kt)/k^2) = \frac{1}{k^2}$
 $\lim_{k \to \infty} ||f_k - 0||_\infty = \lim_{k\to\infty} \frac{k}{k^2} = 0$
So, $f_k \to 0$ as $k \to \infty$.

Unfortunately, not all normed space is not "closed" under limit operation.



We call a complete normed space a Banach space. Examples of Banach spaces:

- IRⁿ with any norm
- · C[a,b] with 11.1100

§ 1.3 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 goad 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 IR" because it is simple yet able to model a variety of data sets (e.g., singals, images, videos, attributes of things) · Actually, the methods can be extended to any normed spaces. · Applications: - Topic discovery. Suppose the Niectors are word histograms with N documents respectively, i.e., the j-th component in X2. is the counters of the j-th word in document i. A clustering algorithm patition the documents into K groups, which typically can be intepreted as groups of documents with the same topics, genre, or author. - Patient Austering. If {Xiliare feature vectors associated with N patients admitted to a hospital, a clustering algorithm clusters the patient into K groups of similar patients. - Recommandation system. A group of N people respond to ratings of n movies. A clustering algorithm can be used to cluster the people into K groups, each with similar taste.

Then we can recommend new movies liked by someone to people in the same group as him/her. - Many other applications. Mathematical formulation: · Representation: Let Ci E {1, 2, -, k} be the group that Xi belongs to . i=1, 2, --, N. Then, group j, denoted by Gij, is $G_j = \{i | C_i = j\}, j = 1, 2, \dots, k$. We assign each group a representative vector, denoted by Z1, Z2,..., ZK. The representative vectors are not necessarily one of given vectors. · Evaluation: First of all, within one specific group j G;, all vectors should be close to the representative vector Zj. More precisely, let $J_j = \sum_{i \in G_j} \|\chi_i - z_j\|_2^2$ Then, J; should be small. Secondly, consider all groups, since each J; is small, $J = J_1 + J_2 + \cdots + J_{\kappa}$ should be small. Altogether, we solve the following $\begin{array}{c} \underset{G_{i},\cdot,G_{k}}{\text{min}} & \underset{G_{j},\cdot,G_{k}}{\overset{K}{\rightarrow}} J_{j} \\ \underset{Z_{1},\cdot,Z_{k}}{\overset{K}{\rightarrow}} & \underset{Z_{i},\cdot,Z_{k}}{\overset{K}{\rightarrow}} J_{j} \end{array} \xrightarrow{} \begin{array}{c} \underset{G_{i},\cdot,G_{k}}{\text{min}} & \underset{J=1}{\overset{K}{\rightarrow}} \left(\underset{i\in G_{j}}{\overset{K}{\rightarrow}} \|\chi_{i} - Z_{j}\|_{2}^{2} \right) \\ \underset{Z_{i},\cdot,Z_{k}}{\overset{K}{\rightarrow}} \end{array}$ · Optimization. We may use an alternating minimization to solve the minimization. Step 1: Fix the representatives 2,..., 2x, find the best partitions GI, .., GK, i.e., solve $\min_{G_{i}, \dots, G_{k}} \sum_{j=1}^{k} \left(\sum_{i \in G_{i}} \|\chi_{i} - Z_{j}\|_{2}^{2} \right) - - - (1)$

$$\begin{array}{c} \text{Step 2: Fix the groups } G_{1}, \cdots, G_{K}, \quad \text{find the best representatives} \\ & & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ &$$

from
$$\chi_1, \chi_2, ..., \chi_N$$
 randomly.
Step 1: Given $z_1, z_2, ..., z_k$, compute
 $C_i = \arg\min_{j \in \{1, 2, ..., N\}} ||\chi_i - z_j||_2^2$, $i=1, 2, ..., N$.
and define
 $G_j = \{i \mid C_i = j\}$, $j=1, 2, ..., k$.
Step 2: Given $G_i, G_2, ..., G_K$, compute
 $z_j = \frac{1}{|G_j|} (\sum_{i \in G_j} \chi_i)$
Go back to Step 1.

K-medians Algorithm
In K-means, the Euclidean norm is used. We can replace it by
(-norm. We solve

$$\begin{array}{c} \min_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \end{array}$$
The numerical solver is
Step 1: Fix $Z_{1,i}, \cdots, Z_{K}$, Solve

$$\begin{array}{c} \min_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{j}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{j}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{i}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{i}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \left(\sum_{i \in G_{i}} ||X_{i} - Z_{j}||_{i} \right) \\ \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}} \sum_{j=1}^{K} \sum_{\substack{G_{1,i},G_{K} \\ i \in G_{i}}}$$

Similar to the discussion in k-means, it is decomposed
into K sub problems
$$\begin{array}{c} \underset{2j}{\min} & \underset{i \in G_j}{\sum} \|[\chi_i - \mathcal{Z}_j]\|_i & j = 1, 2, \cdots, K. \\ \end{array}$$
It is well known (Galileo) that the solution is
$$Z_j = median (\chi_i) \\ \underset{i \in G_j}{\sum} & \underset{i \in G_j}{\max} (\chi_i) \\ \end{array}$$
Where median (χ_i) takes component-wise median.
This algorithm is called "k-median" algorithm.

\$ 1.4 Inner product / Hilbert space
Norms give only metrics, i.e., measuring the distance of two vectors.
However, in many applications, the "angel" of two vectors matter.
· For example; two images X, y showing the same scene with different lights.
For simplicity, we may assume, say, $y = \frac{1}{2} \chi$ (the first image is with 100%)
Then $ \chi-y = \frac{1}{2} \chi \longrightarrow not small.$
but x, y are from the same scene.
· We use inner product for "ange(" of two vectors
٠
Definition; A function $\langle \cdot, \cdot \rangle : V \times V \rightarrow IR$ is called an inner product
on the vector space V over R if

(2) $\langle x_{1} + \beta x_{2}, y \rangle = \alpha \langle x_{1}, y \rangle + \beta \langle x_{2}, y \rangle$ $\forall d, \beta \in \mathbb{R}, x_{1}, x_{2}, y \in V.$ (3) $\langle x, y \rangle = \langle y, x \rangle$

Example 1:
$$\mathbb{R}^n$$
 is a vector space. We can define an inner product as
 $\langle x, y \rangle = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$, where $x = \begin{bmatrix} x_1 \\ x_n \end{bmatrix}$ and $y = \begin{bmatrix} y_1 \\ y_n \end{bmatrix}$.
 $(\equiv x^T y)$

Example 2: Another inner product in
$$\mathbb{R}^n$$
 is as follows.
Let $A \in \mathbb{R}^{n \times n}$ be a symmetric positive definite (spd)
(Recall spd means: $A^T = A$ and $x^T A X > 0 \quad \forall \quad x \neq 0$)
Then $\langle x, y \rangle_A = x^T A y$ defines an inner product in \mathbb{R}^n , keause
 $0 \quad \langle x, x \rangle_A = x^T A X \ge 0$ and $\langle x, x \rangle_A = 0 \Leftrightarrow x^T A x = 0 \Leftrightarrow x = 0$.
(2) $\langle dx_1 + \beta x_2, y \rangle_A = (dx_1 + \beta x_2)^T A y = \alpha \times_1^T A y + \beta \times_2^T A y$
 $= \alpha \langle x_1, y \rangle_A + \beta \langle x_2, y \rangle_A$.
(3) $\langle x, y \rangle_A = x^T A y = (x^T A y)^T = y^T A^T x = y^T A x = \langle y, x \rangle_A$.
Example 3: In $\mathbb{C}[a_1b]$, we can define an inner product as
 $\langle f, g \rangle = \int_{\infty}^{b} f(x) g(x) dx$, $\forall f, g \in C[a_1b]$.

Cauchy -Schwartz Inequlity:
If
$$\langle \cdot, \cdot \rangle$$
 is an inner product on V, then, for any $\chi, y \in V$,
 $|\langle \chi, Y \rangle|^2 \leq \langle \chi, \chi \rangle \langle Y, Y \rangle$.
The equality holds true if and only if $\chi = \alpha Y$ for some $\alpha \in \mathbb{R}$.
Proof. Let $\chi \in \mathbb{R}$ be an arbitrary number
 $0 \leq \langle \chi + \chi Y, \chi + \chi Y \rangle = \langle \chi, \chi \rangle + \chi \langle Y, \chi \rangle + \chi \langle \chi, Y \rangle + \chi^2 \langle Y, Y \rangle$
 $= \langle \chi, \chi \rangle + 2\chi \langle \chi, Y \rangle + \chi^2 \langle Y, Y \rangle$
Thus, $\chi^2 \langle Y, Y \rangle + 2\lambda \langle \chi, Y \rangle + \langle \chi, \chi \rangle \geq 0$.
 $\chi = \langle \chi, \chi \rangle + 2\chi \langle \chi, Y \rangle + \chi^2 \langle Y, Y \rangle$
Thus is a quadratic function of χ and always non-negative.
 $\chi = \chi \langle \chi, Y \rangle^2 - 4 \langle Y, Y \rangle \langle \chi, \chi \rangle \leq 0$
Finally, when $\langle \chi, Y \rangle^2 = \langle \chi, \chi \rangle \langle Y, Y \rangle$, there is a root, i.e.,
 $\exists a unique \lambda \in \mathbb{R}$, $\chi^2 \langle Y, Y \rangle + 2\lambda \langle \chi, Y \rangle + \langle \chi, \chi \rangle = 0$

$$\begin{array}{c} & & \\ \exists a \text{ unique } \lambda \in \mathbb{R}, \quad \langle x + \lambda y , x + \lambda y \rangle = 0 \\ & & \\ & & \\ \exists a \text{ unique } \lambda \in \mathbb{R}, \quad x + \lambda y = 0 \\ & & \\ & \exists a \text{ unique } \lambda \in \mathbb{R}, \quad \chi = -\lambda y \end{array}$$

With the Cauchy-Schwartz inequality, we can show that:

$$\begin{array}{c} \|X\| = \langle x, x \rangle^{k_{2}} & defines \ a \ norm. \end{array} \longrightarrow \begin{array}{c} Called \quad "norm \ induced \\ & by \ the \ inner \ product". \end{array}$$

$$\begin{array}{c} Proof. \quad D \quad \|X\| = (\langle x, x \rangle)^{k_{2}} \ge 0 \quad and \ \|x\| = (\langle x, x \rangle)^{k_{2}} = 0 \iff x = 0. \end{array}$$

$$\begin{array}{c} @ \quad \|\alpha x\| = \langle \alpha x, \alpha x \rangle^{k_{2}} = (\alpha^{2} \langle x, x \rangle)^{k_{2}} = [\alpha| \ \|x\|| \\ \hline @ \quad \|x + y\|_{2}^{2} = \langle x + y, x + y \rangle = \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle \\ & = \|x\|^{2} + \|y\|^{2} + 2\langle x, y \rangle \qquad \left(\begin{array}{c} Note \ that \ Cauchy-Schworts \\ becomes \end{array} \right) \\ & = (\|x\|^{2} + \|y\|^{2} + 2 \ \|x\|| \ \|y\|| \\ & = (\|x\| + \|y\|)^{2} \end{array}$$

Cauchy Schwortz is restated as:

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

define the angels between
$$x$$
 and y .
• When $x = \alpha y$ with $\alpha > 0$,
 $\begin{array}{c} \langle x, y \rangle \\ ||x||| ||y|| = 1 \\ \hline x y \end{array}$
Since x and y are in the same direction,
It is naturally to define $\langle (x,y) = 0$
• When $x = \alpha y$ with $\alpha < 0$
 $\begin{array}{c} \langle x, y \rangle \\ ||x||| ||y|| = -1 \\ \hline x \end{array}$
Since x and y are in the opposite direction,
it is naturally to define $\langle (x,y) = \overline{n}$.
We define $\cos \langle (x,y) = \frac{\langle x, y \rangle}{||x|| ||y||}$
or, equivalently, $\sum (x,y) = \arccos \frac{\langle x, y \rangle}{||x|| ||y||}$
This definition Coincides with the above two cases and the
vectors in \mathbb{R}^2 or \mathbb{R}^3 equipped with the standard inner product.
Orthogonality: Let V be an inner product space.
• When the angel of x and y is $\overline{\pm}$, we call they are
orthogonal, denoted by $x \pm y$, i.e.,
 $[x \pm y]$ if $\langle x, y \rangle = 0$
• When $x \pm y$, they are least relevant.

· Pythagoras' theorem: Let X, y be two vectors in an inner product space.





Hilbert space: A Hilbert space is a Banach space in which the norm is induced by an inner product. Vector spaces normed spaces product space Hilbert spaces Banach spaces Example 1: IR" with inner product $\langle x, y \rangle = \chi^T y$

is a Hilbert space. Example 2: Rⁿ with inner product $\langle x, y \rangle_A = \chi^T A y$, where A is an spot matrix is a Hilbert space. The norm on this space is $\|\chi\|_{A} = (\chi^{T}A\chi)^{k}.$ Example 3: C[a, b] with inner product $\langle f,g \rangle = \int_a^b f(x)g(x)dx$ $\langle f,g \rangle = \int_a f(x)g(x)dx$ is NOT a Hilbert space, because it is not complete (of a convergent sequence may not be in CEa, b] To complete CTa, b] under the norm $||\cdot|| = \langle \cdot, \cdot \rangle^k$, we need to extend the Riemmanian integral to the so-called Lebesgue integral, and the resulting Hilbert space is $L^2(a,b)$. In the following, we will consider calculus on Hilbert/Banach spaces.

§ 1.5 Case Study: Kernel trick, Kernel K-means The k-means will not work for the following examples Recall in a chastering, we want to group X1, X2, ..., XNER" into K groups. The k-means algorithm works like: Initialize Z1, Z2, --, ZK step 1: Given Z1, Z2, --, ZK, update the groups G1,..., GK by D for each Xi, assign (i, the groups that Xi belongs to, by $C_{i} = \arg\min_{j \in \{i, \dots, k\}} ||X_{i} - z_{j}||_{2}^{2}$ @ Then Gi = {i | Ci = j}, for j=1,2,..., K. step 2: Given Gi, ..., Gir, update their representives by $Z_j = \frac{1}{|G_j|} \left(\sum_{i \in G_j} \chi_i \right) = for \quad j = 1, 2, \cdots, K.$ To modify k-means to those "curved" data sets in IR" we use a transform to "un-curve" the data sets in a Hilbert space. Let $\phi: \mathbb{R}^n \to H$ $-\phi(x_i)$ φ(x;

$$\phi(X_i)$$
 is called the feature of X_i .
 ϕ is called the feature map .
[H is called the feature space.

Then we apply k-means to
$$\phi(x_i), \phi(x_2), ..., \phi(x_n)$$
 in H.
Let $Z_1, ..., Z_k$ be the representative vectors in H.
Step 1: Given $Z_1, ..., Z_k$,
 $C_i = \arg\min_{\substack{j \in \{b_2, \cdot, k\}}} ||\phi(x_i) - Z_j||^2, \quad i=1,2,..,N.$
and
 $G_j = \{i \mid C_i = j\}, \quad j=1,2,..,k.$
Step 2: Given $G_{11}, ..., G_k$,
 $Z_j = \frac{1}{|G_j|} (\sum_{i \in G_j} \phi(x_i))$
Repeat.

However, finding the feature map \$\$ is not easy, because of depends on the shape of X1, X2, ~, XN, which generally is very complicated.

The good news is that: [There is no need to know of explicitly in the K-means algorithm.] This is seen as in below: · First of all, since we care only the groups of X1, .- , XN, we only need to know GI, ..., GK. The representatives Z, .-. ZK are only intermediate. Therefore, we can eliminate 8,, -, 2 in the K-means algorithm.

· Now, only D involves the feature mapping of. Since H is a Hilbert space, we can expand the norm in D by $\left\| \phi(x_i) - \frac{1}{|G_i|} \sum_{\substack{i \in G_i}} \phi(x_i) \right\|^2$ $= \langle \phi(\chi_i) - \frac{1}{|G_j|} \underset{\substack{\ell \in G_j}}{\sum} \phi(\chi_\ell), \phi(\chi_i) - \frac{1}{|G_j|} \underset{\substack{\ell \in G_j}}{\sum} \phi(\chi_\ell) \rangle$ $= \langle \phi(x_i), \phi(x_i) \rangle - \frac{2}{|G_1|} \sum_{\ell \in G_1} \langle \phi(x_i), \phi(x_\ell) \rangle$ $+\frac{1}{|G_{j}|^{2}}\sum_{l_{i}\in G_{i}}\langle\phi(\chi_{l_{i}}),\phi(\chi_{l_{i}})\rangle$ We see that Only inner products in the feature space are involved Therefore, An explicit expression of \$\$ NOT necessary.

Kernel trick: Instead of defining $\phi(x)$ explicitly, we define a kernel function K(x,y), which satisfies $K(x,y) = \langle \phi(x), \phi(y) \rangle$. The kernel function K(x,y) can be seen as an quantification of similarity of x and y.

Not all function
$$k(x, y)$$
 satisfying $k(x, y) = \langle \phi(x), \phi(y) \rangle$ for some
feature map ϕ . (Example: $k(x, y) = -1$ is not yod because
it violates with $\langle \phi(x), \phi(y) \rangle \geq 0$)
Which function $k(x, y)$ can be an inner product $\langle \phi(x), \phi(y) \rangle$ for
some feature mapping ϕ ?
• First of all, inner product property
 $K(x, y) = \langle \phi(x), \phi(y) \rangle \leq \langle \phi(y), \phi(x) \rangle = k(y, x)$.
(We say $k(\cdot, \cdot)$ is symmetric if $k(x, y) = k(y, x)$ for all $x, y \in \mathbb{R}^n$).
• Secondly, (at $y, y_2, ..., y_m$ be m vectors in \mathbb{R}^n , then,
for any $C = {i \choose 0} \in \mathbb{R}^m$,
 $\langle \sum_{i=1}^m C_i \phi(y_i), \sum_{i=1}^m C_i \phi(y_i) \rangle \geq 0$ (By inner product property)
Dn the other hand,
 $\langle \sum_{i=1}^m C_i \phi(y_i), \sum_{i=1}^m C_i \phi(y_i) \rangle = \langle \sum_{i=1}^m C_i \phi(y_i), \sum_{j=1}^m C_j \phi(y_j) \rangle$
By inner product $= \sum_{i=1}^m \sum_{j=1}^m C_i G_j \langle f(y_i, y_j) - k(y_i, y_m) - k(y_m, y_m) \rangle$
 $= C^T \begin{bmatrix} k(y_m, y_m), k(y_m, y_m) - k(y_m, y_m) \\ k(y_m, y_m), k(y_m, y_m) - k(y_m, y_m) \end{bmatrix}$
 $C \geq 0$. $\psi \in \mathbb{R}^n$.
We say a function $k(\cdot, \cdot) : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$ is symmetric positive
semi-definite if:
 i $k(x, y) = k(y, x)$ $\forall x, y \in \mathbb{R}^n$
 i $k(x, y) = k(y, x)$ $\forall x, y \in \mathbb{R}^n$
 i $k(x, y_m) = k(y_m, y_m) - k(y_m, y_m) - k(y_m, y_m)$
 $k(y_m, y_m) - k(y_m, y_m) - k(y_m, y_m)$
 i i $k(x, y) = k(y_m) \times i \times x_m + y_m \in \mathbb{R}^n$, the matrix
 $\begin{pmatrix} k(y_m, y_m), k(y_m, y_m) - k(y_m, y_m) -$

Mercer's theorem tells us that: If a kernel function
$$K(\cdot, \cdot)$$
 is
symmetric positive semi-definite, then there exists a feature map
 ϕ such that $K(x,y) = \langle \phi(x), \phi(y) \rangle$

Some popular kernels:
D
$$K(X,y) = X^T y$$
 $(\phi(x) = X$. No transform)
D $K(X,y) = (x^T y)^{\alpha}$ polynomial kernels
O $K(x,y) = e^{-\frac{\|(x-y)\|_{2}^{2}}{\sigma^{2}}}$ Gaussian kernel

Kernel k-means algorithm
• choose a kernel function
$$K(\cdot, \cdot)$$

• Initialize $G_{I}, G_{2}, \cdot \cdot, G_{K}$ by , e.g. one step of k-means.
>• Set
 $C_{i} = \arg\min_{j \in \{h_{2r}, k\}} \left(K(x_{i}, x_{i}) - \frac{2}{|G_{j}|} \sum_{k \in G_{j}} K(x_{i}, x_{k}) + \frac{1}{|G_{j}|^{2}} \sum_{k \in G_{j}} K(x_{k}, x_{k}) \right)$
for $i=1, 2, -\cdot, N$.
• update $G_{I}, G_{2}, -\cdot, G_{K}$ by
 $G_{j} = \{i \mid C_{i}=j\}$, for $j=1, 2, \cdots, K$.
• Go back and repeat

Example: If we use Gaussian kernel $K(x,y) = e^{-\frac{11x-y_1y^2}{\sigma^2}}$ then • $K(X_i, \chi_i) = e^{-\frac{||\chi_i - \chi_i||_2}{\sigma^2}} = 1$ - so all \$(X1), ..., \$(Xw) are on unit sphere in 1-1. $k(x_i, x_j) \begin{cases} \approx 0 & \text{if } ||x_i - x_j||_2 \text{ is large} \\ \approx | & \text{if } ||x_i - x_j||_2 \text{ is small}. \end{cases}$ - so $\phi(x_i)$, $\phi(x_j)$ are orthogonal in 1-1 if $\|X_i - X_j\|_2$ large $-\phi(x_i) \approx \phi(x_j)$ in (+ if $||x_i-x_j||$ small. Therefore, R" 1-1 Thus, K-means works for this data set.