Linear functions and Differentiable functions

for f being linear on V.

Properties of linear functions  
• Homogeneity: 
$$f(\sigma x) = \alpha f(x)$$
  $\forall x \in \mathbb{R}$  and  $x \in V$ .  
(Because  $f(\alpha x) = f(\alpha x + \alpha x) = \alpha f(x) + 0 f(x) = \alpha f(x)$ )  
It implies  $f(0) = 0$ , because  $f(0) = f(0 \cdot x) = 0$   $\forall x \in V$ .  
• Additivity :  $f(x + y) = f(x) + f(y)$   $\forall x, y \in V$ .  
•  $f(\alpha_1 x_1 + \dots + \alpha_k x_k) = \alpha_1 f(x_1) + \dots + \alpha_k f(x_k)$ ,  $\forall \alpha_1, \dots, \alpha_k \in \mathbb{R}$ ,  $x_1, \dots, x_k \in V$ .  
To see this, we note that  
 $f(\alpha_1 x_1 + \dots + \alpha_k x_k) = \alpha_1 f(x_1) + f(\alpha_2 x_2 + \dots + \alpha_k x_k)$   
 $= \alpha_1 f(x_1) + \alpha_2 f(x_2) + f(\alpha_3 x_3 + \dots + \alpha_k x_k)$ 

Innex product representation of a linear function on Hilbert spaces  
For simplicity, let's consider a linear function on 
$$\mathbb{R}^n$$
 equipped with the  
standard innex product  $\langle X, Y \rangle = x^T Y$  and the induced norm  $\|X\|_2 = (\langle X, X \rangle)^{\ell_2}$   
• From the discussion above,  
For any given  $A \in \mathbb{R}^n$ , 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  $A \in \mathbb{R}^n$ .  
To see this, let  $e_1, e_2, \dots, e_n$ , where  $e_i = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \in ith$ , be a basis of  $\mathbb{R}^n$ ,  
So that any  $X = \begin{pmatrix} X_1 \\ X_2 \\ X_N \end{pmatrix} \in \mathbb{R}^n$  is written as  $X = X_1 e_1 + X_2 e_2 + \dots + X_n e_n$ .  
Therefore, if f is a linear function, then

$$f(x) = f(x_{i}e_{i} + x_{i}e_{i} + \dots + x_{n}e_{n})$$

$$= x_{i}f(e_{i}) + x_{2}f(e_{i}) + \dots + x_{n}f(e_{n}) - by property of binary functions = , 
where  $a = \begin{pmatrix} f(e_{i}) \\ \vdots \\ f(e_{n}) \end{pmatrix} \in \mathbb{R}^{n}$ .  
• Furthermore, the representation of a linear function  $f(x) =$  is unique, which means there is only one vector  $a \in \mathbb{R}$  for which  $f(x) =$  holds for all  $x$ .  
Indeed, suppose that  $a$  is not unique, i.e., we have two vectors  $a_{i}b$  such that  $f(x) =$  and  $f(x) =$  for all  $x \in \mathbb{R}^{n}$ .  
Then, let  $x = e_{i}$ :  $f(e_{i}) =  = a_{i}$  and  $f(e_{i}) = = b_{i}$ .  
• Aleogeter, we see that  
 $a$  linear function  $f: \mathbb{R}^{n} \rightarrow \mathbb{R}$  if and only if  $f(x) =$  for some unique  $a \in \mathbb{R}^{n}$ .  
The above holds true for any linear function on Hilbert spaces, widdy known as Riesz representation theorem.  
Theorem (Riesz representation thareom):  
Let H be a Hilbert space, and f be a function:  $H \rightarrow \mathbb{R}$ . Then  
 $f$  is linear and bonduli f and only if  $f(x) =$  for some unique  $a \in H$ .  
Example 1: We know mean(x) is linear on  $\mathbb{R}^{n}$ . Since  $\mathbb{R}^{n}$  is a Hilbert space, we can find a unique  $a \in \mathbb{R}^{n}$  st.  
 $mean(x) = t(x+x,t-...tx_{n}) = tx_{n} + tx_{n} =$$$

where 
$$a = (\pm, \pm, \cdots, \pm)^T$$
.  
Example 2: Let H be a Hilbert space, and II-II is the norm.  
It is known that the norm function is NOT kinear.  
Therefore,  
there obesn't exist  $a \in H$  such that  $||x|| = \langle a, x \rangle$  bixeH.  
Hyperplanes  
• Again, we consider  $\mathbb{R}^n$  as a Hilbert space, where any binear  
function is written as  $\langle a, x \rangle$  for some  $a \in \mathbb{R}^n$ .  
Consider the set  
 $S_{a,p} = \{ x \mid \langle a, x \rangle = 0 \}$ ,  
Then,  $\forall x, y \in S_{a,p}$  and  $\forall, \beta \in \mathbb{R}$ ,  
 $\langle a, \alpha \langle x + \beta \rangle \rangle = \alpha \langle a, x \rangle + \beta \langle a, y \rangle = 0 \Rightarrow \alpha \langle x + \beta \rangle \in S$ .  
Thurefore,  $S_{a,p}$  is a plane.  
Since the co-dimension of So is 1 (because it is define by one equation)  
 $S_{a,p}$  is called a hyperplane.  
Now lets consider  
 $S_{a,b} = \{ x \mid \langle a, x \rangle = b \}$  for  $b \in \mathbb{R}$  is given.  
Let  $x_0 \in S_{a,b}$ , i.e.,  $\langle a, x_0 \rangle = b$ , be fixed.  
Then  $S_{a,b} = S_{a,0} + x_0$  because:  
(D)  $\forall x \in S_{a,b} = \langle a, x \rangle = \langle a, x \rangle - \langle a, x_0 \rangle = 0$ ,  $\Rightarrow x + x_0 \in S_{a,0}$ .

 $(2) \forall x \in S_{a,o} \quad \langle a, x + \chi_o \rangle = \langle a, x \rangle + \langle a, \chi_o \rangle = b \implies x + \chi_o \in S_{a,b}$ In other words, Sa, b is a shift of a hyperplane, still called a hyperplane. Sab Xo 0 Sa,0 • This concept can be generalized to any inner product space V. The set {XEV < a, x>=b}, where a EV and bER are given is called a Hyperplane in V. Projection onto hyperplanes · Consider a Hilbert space V and a hyperplane S  $S = \{x \in V \mid \langle a, x \rangle = b\}$ Let YEV be a given vector. The vector on Sthat is the closest to Y is called the projection of yon S, denoted by Psy, Psy i.e.,  $P_s y = \arg \min_{x \in s} ||x - y||$ . · Let us find an explicit expression of Psy in terms of a, b, and y. Theorem: Z is a solution of min ||x-y|| if and only if ZES and <Z-Y, Z-x>=0 VXES. Proof DWe first prove that: If ZES is a solution of min 11x-y11,

then <Z-Y,Z-X>=0 YXES.

Since 
$$\Xi$$
 is a solution,  $\Xi \in S$ , i.e.,  $\langle a, Z \rangle = b$ .  
 $\forall x \in S$  and  $t \in R$ , it is easy to see that  
 $\langle a, (|+t) \ge -tx \rangle = (i+t) \langle a, Z \rangle - t \langle a, x \rangle = b$ .  
Therefore,  $(i+t) \ge -tx \in S$ .  
Since  $Z$  is closest to  $Y$  on  $S$ , we have  
 $\|II \ge -Y\|^2 \le \|((i+t) \ge -tx - Y)\|^2$   
 $= \|(Z - Y) + t(Z - x)\|^2$   
 $= \|(Z - Y) + t(Z - x)\|^2$   
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t^2}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < Z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .  
i.e.,  $t < z - Y$ ,  $Z - x \rangle \ge -\frac{t}{2} \|Z - x\|^2$ .

Ø

Theorem: The solution of $\min_{x \in S}   x - y  ^2$ exists and unique, which is
given by $y = \left(\frac{\langle a, y \rangle - b}{  a  ^2}\right)a$
proof. denote $Z = y - (\frac{\langle a, y \rangle - b}{a})a$ .
$D \langle a, z \rangle = \langle a, y \rangle - \left( \frac{\langle a, y \rangle - b}{  a  ^2} \right) \langle a, a \rangle$
= (a,y> - ((a,y7-6) = 6, so ZES
② ∀ x∈S,
$(z-y, z-x) = -\frac{\langle a, y-b}{  a  ^2} \langle a, z-x \rangle$
$= -\frac{\langle a, y \rangle - b}{  a  ^2} \left( \langle a, z \rangle - \langle a, x \rangle \right) = 0$
$(because \langle a, z \rangle = \langle a, x \rangle = b)$
By the previous theorem, $Z$ is a solution of $\max_{X \in S}   X - Y  $ .
It remains to show the uniqueness.
Suppose we have two solutions $z_1$ and $z_2$ . Then,
$Z_1$ is a solution, $= X Z_1 - Y$ , $Z_1 - Z_2 > = 0$
$Z_2$ is a solution, $\Rightarrow \langle Z_2 - Y, Z_2 - Z_1 \rangle = 0$
Taking difference leads to $(z_1 - z_2, z_1 - z_2) = 0$
$\implies   \mathcal{Z}_1 - \mathcal{Z}_2  ^2 = 0 \implies \mathcal{Z}_1 = \mathcal{Z}_2 \qquad \bigotimes$
• In summary, the projection Psy of YEV onto the hyperplane
$S = \{x \in V   \langle a, x \rangle = b\}$
exists and is unique. Furthermore,
$P_{s}y = y - (\frac{\langle a, y \rangle - b}{  a  ^2})a$
and it satisfies
$\langle P_s y - y, P_s y - x \rangle = 0.$



Affine functions  
A linear function plus a constant is called an affine function.  
That is, a function 
$$f: V \Rightarrow R$$
 is affine if  
 $f(x) = g(x) + b$ ,  
where  $g: V \Rightarrow R$  is linear and  $b \in R$  is a anstant.  
Properties:  
If  $f: V \Rightarrow R$  is affine, then  
 $f(xx + \beta y) = \alpha f(x) + \beta f(y)$   $\forall x, y \in V$  and  $\alpha, \beta \in R$  s.t.  $\alpha + \beta = 1$ .  
To see this, linear  $\alpha + \beta = 1$   
 $f(\alpha x + \beta y) = g(\alpha x + \beta y) + b = \alpha g(\alpha) + \beta g(y) + (\alpha + \beta) b$   
 $= \alpha (g(\alpha) + b) + \beta (g(y) + b) = \alpha f(x) + \beta F(y)$ .  
If  $f: V \Rightarrow R$ , where  $V$  is a Hilbert space, then  
 $f$  must be in the form of  
 $f(x) = \langle a, x \rangle + b$ , where  $a \in V$  and  $b \in R$ .

§ 2.2 Case Studies: Regression and Classification.

§2.2.1 Regression: · Given a set of data  $(\chi_1, y_1), (\chi_2, y_2), \dots, (\chi_N, y_N),$ where Xi EIR<sup>n</sup> is an input feature vector , i=1,2,...,N. YiER is the corresponding response to Xi. Given a new input feature vector XER", how to predict the corresponding response YER? For example, Xi ER represents n attributes of a house, and YiER is the selling price. We want to predict the selling price of a house with feature XER". · Mathematically, we need to find a function  $f:\mathbb{R}^n \to \mathbb{R}$  such that  $f(X_i) \approx Y_i$  ,  $i=1,2,\cdots,N$ This is called regression. In this context, Xi, are called regressor / independent varibles Yi are called dependent variables / out come / label. • The class of all functions  $\mathbb{R}^n \rightarrow \mathbb{R}$  is too large, and the given data set {(Xi, Yi)}\_{i=1}^N is not enough to determine a function uniquely. So, we need to find a function class  $\Phi$  where we search f. Intuitively, larger N, larger function class  $\overline{\Phi}$ . · Linear model: We search f in the class orf all affine functions, i.e.,  $f(x) = \langle a, x \rangle + b$  for some  $a \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ .

Thus, we find 
$$\alpha \in \mathbb{R}^{n}$$
 and  $b \in \mathbb{R}$ , s.t.  
 $\langle \alpha, \chi_{i} \rangle + b \approx y_{i}$ ,  $i=1,2,\cdots,N$ ,  
by minimizing the error of the linear equations.  
While there are many possible definitions of error, it is popular to  
consider the square error as follows:  
 $(\langle \alpha, \chi_{i} \rangle + b - y_{i})^{2}$ ,  $i=1,2,\cdots,N$ .  
Therefore, we find  $\alpha \in \mathbb{R}^{n}$ ,  $b \in \mathbb{R}$  by solving  
 $\min_{\substack{\alpha \in \mathbb{R}^{n}, i} \in \mathbb{R}^{n}$ ,  $b \in \mathbb{R}$  by solving  
 $\min_{\substack{\alpha \in \mathbb{R}^{n}, i} \in \mathbb{R}^{n}$ ,  $b \in \mathbb{R}^{n}$ ,  $b \in \mathbb{R}^{n+1}$   
 $\lim_{\substack{\alpha \in \mathbb{R}^{n}, i} i \in \mathbb{R}^{n \times (n+1)}$ ,  $\beta = \begin{bmatrix} \alpha \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$   
 $\lim_{\substack{\alpha \in \mathbb{R}^{n}, i} i \in \mathbb{R}^{n}$ .  
This problem is called the least squares ( $LS$ ) problem.  
Write  $\chi = \begin{bmatrix} \chi_{i}^{T} & i \\ \chi_{i}^{T} & i \end{bmatrix} \in \mathbb{R}^{N \times (n+1)}$ ,  $\beta = \begin{bmatrix} \alpha \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$   
and  $y = \begin{pmatrix} y_{i} \\ y_{i} \end{pmatrix} \in \mathbb{R}^{N}$ .  
Then  $LS$  problem becomes  
 $\begin{bmatrix} \min_{\substack{\beta \in \mathbb{R}^{n+1} \\ \beta \in \mathbb{R}^{n}$ .  
Then  $LS$  problem becomes  
 $\begin{bmatrix} \min_{\substack{\beta \in \mathbb{R}^{n}, \\ \beta \in \mathbb{R}^{n}} \\ \beta = n+1 \\ \beta \in \mathbb{R}^{n} \\ \beta = n+1 \\ \beta =$ 

Here 
$$\|p\|_{\alpha}^{\alpha}$$
 is the regularization term  
 $\lambda > 0$  is a predefinal regularization parameter.  
In other words, we find  $\beta$  such that frame in  
the error of data fitting and the 2-norm of  $\beta$  regression  
are minimized simultaneously.  
Therefore, ridge regression gives a  $\beta$  such that  
 $X \beta \approx y$  and  $\|\beta\|_{\alpha}$  is small.  
 $- LASSO$  regression: we solve  
 $\begin{bmatrix} minimized \\ pare \\ \\ par$ 

Feature map 
$$\oint: \mathbb{R}^n \rightarrow H$$
  
Then do regression in  $H$ .  
However, Since  $H$  is very large, the set of all linear functions  
is also two large. We need regularization.  
We solve  

$$\frac{\min_{a \in H} \frac{1}{2} \prod_{i=1}^{N} (\langle a, \phi(x_i) \rangle - y_i)^2 + \lambda \|a\|_{H}^{2}}{a \in H}$$
Representer Theorem:  
The soluction must be in the form of  $a = \sum_{i=1}^{N} C_i \phi(x_i)$  for  
some  $C = \begin{bmatrix} C_i \\ C_i \end{bmatrix} \in \mathbb{R}^N$ .  
Proof. For any  $a \in H$ , we claim that  $a$  can be decomposed as  
 $a = a_s + \sum_{i=1}^{N} C_i \phi(x_i)$   
where  $C = \begin{bmatrix} C_i \\ C_i \end{bmatrix} \in \mathbb{R}^N$  and  $\langle a_s, \phi(x_i) \rangle = 0$  for  $i=1,2,...,N$ .  

$$\frac{a_s}{\sum_{i=1}^{N-\ldots-n}}$$
(This is  $N$  linear equation with  $N$  unknows, and it can be  
checked there exists at least one solution)  
Denote  $a_s = a - \sum_{i=1}^{N} C_i \phi(x_i)$ . It can be checked  
 $\langle a_s, \phi(x_i) \rangle = 0$ ,  $j=1,2,...,N$   
due to the construction of  $C$  and  $a_s$ .  
Therefore,  
 $\frac{1}{2} \sum_{i=1}^{N} \langle \langle a, \phi(x_i) \rangle - y_i \rangle^2 + \lambda \|a\|_{H}^{2}$ 

$$= \frac{1}{2} \sum_{i=1}^{\infty} \left( \left\{ \sum_{j=1}^{N} G_{ij}^{2}(x_{j}) + a_{ij}^{2} g(x_{i}) \right\} - y_{i}^{2} + \lambda \right\| \left\{ \sum_{i=1}^{N} G_{ij}^{2}(x_{i}) + a_{ij}^{2} g(x_{i}) \right\} + 2 \left\{ a_{ij}^{2} \sum_{i=1}^{N} G_{ij}^{2}(x_{i}) \right\} + \lambda \left[ \left\{ \sum_{i=1}^{N} G_{ij}^{2} G(x_{i}) \right\} + 2 \left\{ a_{ij}^{2} \sum_{i=1}^{N} G_{ij}^{2}(x_{i}) \right\} + \lambda \left[ \left\{ \sum_{i=1}^{N} G_{ij}^{2} G(x_{i}) \right\} + \lambda \left[ x_{ij}^{2} G(x_{i}) \right\} + \lambda \left[ x_{ij}^{2} G(x_{i}) \right\} + \lambda \left[ x_{ij}^{2} G(x_{i}) \right] + \lambda \left[ x_{ij$$

Let 
$$X+6S+$$
 and  $X-6S_-$  such that  
 $\|[X_{4}-X_{-}\|]_{2} = dist (S_{4}, S_{-}).$   
Since  $X_{+}$  is a projection of  $X_{-}$  onto  $S_{+} = \{x_{1} < a, x\} = 1-b\}$   
 $X_{+} = X_{-} - \frac{\langle a, X_{-} > tb^{-1} \rangle a}{\|a\|_{2}^{2}} a$   
 $= X_{-} - \frac{\langle -1-btb^{-1} \rangle a}{\|a\|_{2}^{2}} a$  (since  $X_{-} \in S_{-}$ )  
 $= X_{-} + \frac{2}{\|A\|_{2}} \alpha$   
Thus,  $\|[X_{+}-X_{-}\|]_{2} = \|\frac{2}{\|A\|_{2}} \alpha$   
Thus,  $\|[X_{+}-X_{-}\|]_{2} = \|\frac{2}{\|A\|_{2}} \alpha$   
Support Vector Machine (SVM)  
 $\max_{x \in \mathbb{N}^{n}} \frac{2}{\pi \|a\|_{x}}$   
 $\xi \in \mathbb{R}$   $\frac{2}{\pi \|a\|_{x}}$   
 $\xi \in \mathbb{R}$   $\frac{2}{\pi \|a\|_{x}}$   
 $\xi \in \mathbb{R}$   $\frac{1}{\pi \|a\|_{x}}$   
 $St. \langle a, X_{1} > tb \geq 1 \text{ if } Y_{1} = 1$   
 $\langle a, X_{1} > tb \leq -1 \text{ if } Y_{1} = -1$ ,  
which is equivalent to  
 $\left[\max_{x \in \mathbb{N}^{n}} \frac{1}{2} \|a\|_{x}^{2}$  (SVM-1)  
 $St. \quad Y_{1}(\langle a, X_{1} \rangle + b) \geq 1$   
The above SVM is NOT robust to noise.  
For example shown on the right,  
even we have only two noisy points, there is  $\sum_{x \in \mathbb{N}^{n}} \frac{\langle x \times x_{-} \times x_{-}$ 

$$\begin{array}{c|c} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\$$

Aagain, one can prove the following representer theorem.  
Theorem: Any solution of 
$$(K-SVM)$$
 is in the form of  
 $a = \sum_{i=1}^{\infty} C_i \phi(X_i)$   
proof. Write  $a = \sum_{i=1}^{\infty} C_i \phi(X_i) + a_s$  for some  $a_s \in H$  and  $\langle a_s, \phi(X_i) \rangle = 0$   
the rest is the same as the linear regression case. We  
Thus,  $(K-SVM)$  becomes  
 $\left[ \min_{C \in \mathbb{R}^N} h\left( y_i \left( \sum_{j=1}^{\infty} K(X_i, X_j) C_j \right) - 1 \right) + \frac{1}{2} C^T K C_r \right)$   
where  $K = [K(X_i, X_j)]_{i=1, j=1}^{N, N} \in \mathbb{R}^{N \times N}$ .  
The prediction of the input  $\chi$  is given by  
 $Sgn\left( \sum_{j=1}^{\infty} K(X, X_j) C_j \right)$   
Again, only  $K(\cdot, \cdot)$  is needed in the kernel SVM,  
and no explicit feature map  $\phi(\cdot)$  is required.

Recall that for a function 
$$f: \mathbb{R} \to \mathbb{R}$$
, the derivative at  $X_0$  is  
 $f'(X_0) = \lim_{x \to X_0} \frac{f(x) - f(X_0)}{x - X_0}$ ,  
which is the same as  
 $\lim_{x \to X_0} \left| \frac{f(x) - f(x_0) - f'(x_0)(x - X_0)}{x - X_0} \right| = 0$   
Notice that  $f(X_0) + f'(X_0)(x - X_0)$  is an affine function in  $\mathbb{R}$  that  
passes through  $(X_0, f(X_0))$ .  
 $f(X_0) + f'(X_0)(x - X_0)$  is an affine function in  $\mathbb{R}$  that  
 $passes$  through  $(X_0, f(X_0))$ .  
 $f(X_0) + f'(X_0)(x - X_0)$   
 $f(X_0) + f'(X_0)(x -$ 

Consider the differentiation of 
$$f$$
 at  $\chi^{o} \in V$ .

(1) By Riesz representation theorem, any affine function is in the form

of 
$$\langle U, x \rangle + a$$
 for some  $U \in V$  and  $a \in \mathbb{R}$ . Since it passes thrue  
 $(\chi^{(0)}, f(\chi^{(0)}), \langle U, \chi^{(0)} \rangle + a = f(\chi^{(0)})$  Therefore, the affine function  
is in the form of  
 $\langle V, \chi \rangle + a = \langle U, \chi - \chi^{(0)} \rangle + \langle U, \chi^{(0)} \rangle + a \rangle$   
 $= f(\chi^{(0)}) + \langle U, \chi - \chi^{(0)} \rangle$ .  
(2). The approxition error is  
 $error = |f(\chi) - f(\chi^{(0)}) - \langle U, \chi - \chi^{(0)} \rangle|_{A}$ 

The error should be in the order of 
$$o(||x-x^{(0)}||)$$
, i.e.,  

$$\frac{error}{||x-x^{(0)}||} \rightarrow o \quad as \quad x \rightarrow x^{(0)}$$

Definition: Let V be a Hilbert space. Let 
$$f: V \rightarrow R$$
. Then  $f$  is  
said Frechet differentiable if there exists a  $v \in V$  such that  
 $\lim_{x \to \chi^{(0)}} \frac{|f(x) - f(\chi^{(0)}) - \langle v, x - \chi^{(0)} \rangle|}{||x - \chi^{(0)}||} = 0$ .  
If  $f$  is differentiable at  $\chi^{(0)}$ ,  $v$  is called the gradient of  $f$   
at  $\chi^{(0)}$ , denoted by  $\nabla f(\chi^{(0)})$ .

Example 1: 
$$f(x) = ||x||^2$$
, where  $||x||$  is the norm on V.  
At any  $\chi^{(o)} \in V$ ,  
 $||x||^2 = ||\langle x - \chi^{(o)} + \chi^{(o)}||^2 = \langle (x - \chi^{(o)}) + \chi^{(o)} \rangle \langle (x - \chi^{(o)}) + \chi^{(o)} \rangle$   
 $= ||x - \chi^{(o)}||^2 + 1|\chi^{(o)}||^2 + 2 \langle \chi^{(o)}, \chi - \chi^{(o)} \rangle$   
Therefore,  $||x||^2 - (||\chi^{(o)}||^2 + 2 \langle \chi^{(o)}, \chi - \chi^{(o)} \rangle) = ||\chi - \chi^{(o)}||^2$   
 $affine approximation$   
So  $\lim_{\chi \to \chi^{(o)}} \frac{||x||^2 - ||\chi^{(o)}||^2 - \langle 2\chi^{(o)}, \chi - \chi^{(o)} \rangle|}{||\chi - \chi^{(o)}||} = \lim_{\chi \to \chi_0} \frac{||x - \chi^{(o)}||^k}{||\chi - \chi^{(o)}||} = 0$ .

Thus, 
$$\nabla f(x^{(0)}) = 2x^{(0)}$$
  
Example 2:  $f(x) = \langle a, x \rangle$  for some  $a \in V$ .  
At any  $x^{(0)} \in V$ ,  
 $\langle a, x \rangle = \langle a, x^{(0)} \rangle + \langle a, x - x^{(0)} \rangle$   
Therefore,  $\lim_{X \to x^{(0)}} \frac{|\langle a, x \rangle - \langle a, x^{(0)} \rangle - \langle a, x - x^{(0)} \rangle|}{||x - x^{(0)}||} = \lim_{X \to 3} \frac{0}{||x - x^{(0)}||} = 0$ .  
Thus,  $\nabla f(x^{(0)}) = a$ .

Example 3: 
$$f(x) = ||x - \alpha||^2$$
 for some  $\alpha \in V$ .  
At any  $\chi^{(o)} \in V$ ,  
 $f(x) = ||x - \alpha||^2 = ||\alpha^{(o)} - \alpha| + (x - x^{(o)})||^2$   
 $= ||x^{(o)} - \alpha||^2 + ||x - x^{(o)}||^2 + 2 \langle x^{(o)} - \alpha, x - x^{(o)} \rangle$   
 $= f(x^{(o)}) + \langle 2(x^{(o)} - \alpha), x - x^{(o)} \rangle + ||x - x^{(o)}||^2$   
So,  $\lim_{x \to x^{(o)}} \frac{|f(x) - f(x^{(o)}) - \langle 2(x^{(o)} - \alpha), x - x^{(o)} \rangle|}{||x - x^{(o)}||} = \lim_{x \to x^{(o)}} ||x - x^{(o)}|| = 0$ .  
Therefore,  $\nabla f(x^{(o)}) = 2(x^{(o)} - \alpha)$ 

Properties:  
(1) Frechet differentiation is linearly, i.e.,  

$$\nabla(\alpha f + \beta g)(x) = \alpha \nabla f(x) + \beta \nabla g(x)$$
.  
(2) Chain rule: Let  $f:V \rightarrow \mathbb{R}$  and  $g:\mathbb{R} \rightarrow \mathbb{R}$ . Then  $gof:V \rightarrow \mathbb{R}$  and  
 $\nabla(gof)(x) = g'(f(x)) \cdot \nabla f(x)$   
if f and g are differentiable at x and  $f(x)$  respectively.  
Example 4:  $f(x) = \Pi X \Pi$   $\forall x \in V$ .

This is a composition of 
$$f_{1}(x) = ||x||^{2}$$
 from  $V \to \mathbb{R}$   
and  $f_{2}(t) = \sqrt{f}$  from  $\mathbb{R} \to \mathbb{R}$ .  
When  $||x|| \neq 0$ , both  $f_{1}$  and  $f_{2}$  are differentiable.  
Also,  $\forall f_{1}(x) = 2\chi$ ,  $f_{2}(t) = \frac{1}{2\sqrt{f}}$  if  $t \neq 0$ .  
So,  $\nabla f(x) = \nabla (f_{2} \circ f_{1})(x) = f_{1}'(f_{1}(x) \cdot \nabla f_{1}(x))$   
 $= \frac{1}{2\sqrt{||x||^{2}}} \cdot 2\chi = \frac{\chi}{||\chi||}$ .  
When  $||\chi|| = 0$ , (i.e.,  $\chi = 0$ ),  $f_{2}(t)$  is NOT differentiable at  $f_{1}(x) = 0$ .  
It can be shown that  $f(x) = ||\chi||$  is NOT differentiable at  $\chi = 0$ .  
(3) For functions on  $\mathbb{R}^{n}$ :  $f : \mathbb{R}^{n} \to \mathbb{R}$   
 $\nabla f(\chi) = \begin{pmatrix} \frac{f_{1}'(\chi)}{2\chi} \\ \frac{f_{2}'(\chi)}{2\chi} \\ \frac{f_{2$ 

§ 2.4. Linear operators and Higher-order darivatives  
• Linear operator / Linear transformations  
Now we consider functions between vector spaces.  
Let V1, V2 be two vector spaces  
A map L: 
$$V_1 \rightarrow V_2$$
 is a linear operator (a.k.a. linear transformation)  
if:  $L(\alpha X + \beta Y) = \alpha L(x) + \beta L(y) \quad \forall \alpha, \beta \in \mathbb{R}, x \in V_1.$   
Example 1: Let  $A = \begin{bmatrix} a_{11} - \cdots & a_{1n} \\ \vdots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$   
Then the transformation:  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  defined by  
 $x \rightarrow AX$ , where  $Ax$  is the matrix-vector product.  
is a linear transformation, because  
 $A(\alpha X + \beta Y) = \alpha AX + \beta AY.$   
Reversely, any linear transformation  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  must be in  
the form of  $L(x) = AX$  for some matrix  $A \in \mathbb{R}^{m \times n}$ .  
Example 2: Let  $f: V \rightarrow \mathbb{R}$  be a linear function on  $V$ .  
Then f is a linear transformation from  $V$  to  $\mathbb{R}$ , as  $\mathbb{R}$  is  
a vector space.  
Example 3: Let  $A \in \mathbb{R}$ . Then define  $L : \mathbb{R} \rightarrow V$  by  
 $L(x) = aX$ . Then  $L$  is a linear operator.  
Example 4: Let  $V_1 = \{f \mid f \text{ ord } f' \text{ is continuous on } Ca, big\}$   
and  $V_2 = \{f \mid f \text{ is an linear operator.} D: V_1 \rightarrow V_2$  defined by  
 $Df = f' \text{ is an linear operator.}$ 

• Operator norm

Consider the set of all linear operators  $V_1 \rightarrow V_2$ , where  $V_1$ ,  $V_2$ are two normed space.

• 
$$\forall A, B$$
 linear operators  $V_1 \Rightarrow V_2$ , define  $A+B$  by  
 $(A+B)(x) = A(x) + B(x)$   $\forall x \in V_1$   
•  $\forall d \in \mathbb{R}$  and A linear operator  $V_1 \Rightarrow V_2$ , define  $A A$  by  
 $(A A)(x) = A A(x)$   $\forall x \in V_1$ .  
Then, the set of all linear operator  $V_1 \Rightarrow V_2$  is a vector space.  
So, we can define a norm on it. For any linear operator  $A: V_1 \Rightarrow V_2$   
 $\|A\| = \sup_{\|x\|_{V_1}=1} \|Ax\|_{V_2}$   $(11 \cdot \|V_1, \|1 \cdot \|V_2 - the norm$   
 $n V_1$  and  $V_2$  respectively  
the unit hall in  $\|1 \cdot \|V_1$ .  
 $V_1$   $V_2$ 



 $\langle A \chi, y \rangle_{V_2} = \langle \chi, A^* y \rangle_{V_1}$  &  $\chi \in V_1$  and  $y \in V_2$ .

2 . ,

Example 1: Consider 
$$A \in \mathbb{R}^{m \times n} = \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$$
  
Then  $\langle A \chi, y \rangle = y^T A \chi = (A^T y)^T \chi = \langle \chi, A^T y \rangle$   $\forall \chi \in \mathbb{R}^n$   
Therefore, the adjoint of  $A$  is its transpose  $A^T$ .  
Example 2: Let  $V$  be an inner product space.  
Consider  $\mathcal{L}(V, \mathbb{R})$ , i.e., bounded linear functions on  $V$ .  
Let  $f: V \Rightarrow \mathbb{R}$ . Then,  $f(\chi) = \langle a, \chi \rangle_V$  for some  $a \in V$ .  
Therefore,  $\forall \chi \in V$ ,  $y \in \mathbb{R}$ ,  
 $\langle f(\chi), y \rangle_{\mathbb{R}} = yf(\chi) = y \langle a, \chi \rangle_V = y \langle \chi, a \rangle_V = \langle \chi, y a \rangle_V$   
Therefore,  $f^*(y) = ya$ ,  $\forall y \in \mathbb{R}$ ,  
Obviously,  $f^* \in \mathcal{L}(\mathbb{R}, V)$ .

• Differentiation of transformations between normed spaces  
Let V1, V2 be two normed spaces with inner products ||.||v, and ||.||v2  
Let F: V1 
$$\rightarrow$$
 V2 be a map (not necessarily linear)  
Then, ct any point  $\chi^{(0)} \in V_1$ , then linear approximation passing thru  
 $(\chi^{(0)}, F(\chi^{(0)})$  is  
 $F(\chi) \approx F(\chi^{(0)}) + L(\chi - \chi^{(0)})$ , where  $L \in L(V_1, V_2)$ .  
If this approximation is  $o(||\chi - \chi^{(0)}||_{V_1})$ , then L is called the differentiation  
Definition: F: V1  $\rightarrow$  V2 is differentiable at  $\chi^{(0)} \in V_1$ , if there exists  
 $L \in L(V_1, V_2)$  such that  
 $\lim_{X \to \chi^{(0)}} \frac{||F(\chi) - F(\chi^{(0)}) - L(\chi - \chi^{(0)})||_{V_2}}{||\chi - \chi^{(0)}||_{V_1}} = 0$ .  
L is called a derivative of F at  $\chi^{(0)}$ , denoted by  
 $DF(\chi^{(0)}) = L$ .

Example 1: If 
$$f: V \rightarrow R$$
 with V a Hilbert space, then  
 $Df(X^{(0)}) y = \langle \nabla f(X^{(0)}), y \rangle \quad \forall y \in V.$   
Example 2: Let  $A \in \mathcal{L}(V_1, V_2)$ , then, for any  $X^{(0)} \in U_1$ ,  
 $\lim_{X \to X^{(0)}} \frac{||AX - AX^{(0)} - A(X - X^{(0)})||_{V_2}}{||X - X^{(0)}||_{V_1}} = 0.$   
 $X \rightarrow X^{(0)} \frac{||X - X^{(0)}||_{V_1}}{||X - X^{(0)}||_{V_1}} = A.$   
Chain Rule: Let  $F: V_1 \rightarrow V_2$ ,  $G: V_2 \rightarrow V_3$ ,  
Then  $G \circ F : V_1 \rightarrow V_3$   
Then  $D(G \circ F)(X) = DG(F(X) \circ DF(X).$   
Example 3:  $f(X) = f_1(X) f_2(X)$ , where  $f, f, f_1: V \rightarrow R$   
Define  $F: V \rightarrow R^2$  as  $F(X) = (f_1(X))$   $\forall X \in V$   
 $G: R^2 \rightarrow R$  as  $G(\mathcal{C}) = \alpha \beta$   $\forall \alpha \beta^{\beta R}$   
Then  $f(X) = G(F(X))$   
So  $Df(X) = DG(F(X)) \circ DF(X)$   
Let's calculate  
 $O DG(\mathcal{C}) = DG(F(X)) \circ DF(X)$   
 $Let's calculate$   
 $O DG(\mathcal{C}) = TG(\mathcal{C}) + Df_1(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $f_1(X) = f_1(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $f_2(X) = f_2(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $V = F(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $V = F(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $V = F(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $V = F(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$   
 $V = F(X^{(0)}) + Df_2(X^{(0)})(X - X^{(0)}) + o(HX - X^{(0)}|)$ 

$$\begin{split} Df(x)(y) &= \left\langle \begin{pmatrix} f_{1}(x) \\ f_{1}(x) \end{pmatrix}, \begin{array}{l} Df_{1}(x)(y) \\ Pf_{2}(x)(y) \\ \end{array} \right\rangle \\ \hline Df(x)(y) &= Df_{1}(x)(y), f_{2}(x) + f_{1}(x), Df_{2}(x)(y) \\ \hline Ix perturbar, if V is a Hilbert space, \\ \hline Vf(x) &= f_{2}(x), \nabla f(x) + f_{1}(x), \nabla f_{2}(x) \\ \hline Vf(x) &= f_{2}(x), \nabla f(x) + f_{1}(x), \nabla f_{2}(x) \\ \hline Vf(x) &= f_{2}(x), \nabla f(x) + f_{1}(x), \nabla f_{2}(x) \\ \hline Vf(x) &= f_{2}(x), \nabla f(x) + f_{1}(x), \nabla f_{2}(x) \\ \hline Let \ us \ calculate \ Df(x^{(u)}) \ and \ \nabla f(x^{(u)}) \\ f(x) &= \frac{1}{2} \|Ax - b\|_{w}^{2} = \frac{1}{2} \|A(x^{(u)} - b) + A(x - x^{(u)})\|_{w}^{2} \\ &= \frac{1}{2} \|Ax^{(u)} - b\|_{w}^{2} + \langle Ax^{(u)} - b, A(x - x^{(u)})\|_{w}^{2} \\ &= f(x^{(u)}) + \left\langle A^{*}(Ax^{(u)} - b), x - x^{(u)} \right\rangle, \\ &+ \frac{1}{2} \|A(x - x^{(u)})\|_{w}^{2} \\ &= f(x^{(u)}) + \left\langle A^{*}(Ax^{(u)} - b), x - x^{(u)} \right\rangle, \\ &= \frac{1}{2} \|A(x - x^{(u)})\|_{w}^{2} \\ &= \frac{1}{$$

• Hessian of functions on Hilbert spaces  
Let 
$$f: V \rightarrow \mathbb{R}$$
 where V is a Hilbert space.  
I-st order derivative :  $\nabla f(x) \in V$ .

2-nd order derivative: view 
$$X \rightarrow \nabla f(x)$$
 as a map  $V \rightarrow V$ , and  
its derivative  $D(\nabla f)$  is the 2nd order derivative of of  $f: V \rightarrow R$ .  
Definition:  $\nabla^2 f(x) \equiv D(\nabla f(x))$   
Therefore,  $\nabla^2 f(x)$  is a bounded linear transformation  $V \rightarrow V$ .  
Example :  $f(x) = \frac{1}{2} ||AX - b||_{W}^2$ , where A bounded linear  $V \rightarrow W$   
 $b \in W$ ,  $x \in V$ .  
Then,  $\nabla f(x) = A^*(AX - b) = A^*AX - A^*b$   
and  $D(\nabla f(x)) = A^*A$   
So,  $\nabla^2 f(x) = A^*A \in \mathcal{L}(V, V)$   
In particular, if  $X \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{m \times n}$   
 $f(x) = \frac{1}{2} ||AX - b||_2^2$   
Then  $\nabla f(x) = A^T(AX - b)$   
 $\nabla^2 f(x) = A^TA$ 

• Function expansion  
• 
$$\frac{d}{dt} f(x+tu) = \langle \nabla f(x+tu), u \rangle$$
, where  $t \in \mathbb{R}$ ,  $x, u \in V$ .  
because  $f(x+su) = f(x+tu) + \langle \nabla f(x+tu), (s-t)u \rangle + o(1s-t)|uu|)$   
 $= f(x+tu) + (s-t) \langle \nabla f(x+tu), u \rangle + o(1s-t)|uu|)$   
Set  $t=0$ :  $\frac{d}{dt} f(x+tu)|_{t=0} = \langle \nabla f(x), u \rangle$   
i.e.,  $\left[ \langle \nabla f(x), u \rangle \text{ is the directional derivative of } f(x) \text{ along } U \right]$   
• Similarly,  
 $\frac{d^2}{dt^2} f(x+tu) = \frac{d}{dt} \langle \nabla f(x+tu), u \rangle = \langle \nabla^2 f(x+tu) u, u \rangle$   
because  $\langle \nabla f(x+su), u \rangle = \langle \nabla f(x+tu), u \rangle$   
 $+ \langle \nabla^2 f(x+tu) (s-t)u, u \rangle + o(1s-t)||uu||^2 \rangle$ 

$$= \langle \nabla f(x+tu), u \rangle + (s-t) \langle \nabla^{2} f(x+tu), u, u \rangle + o(ts-t) \rangle$$
Set  $t=o; \frac{d^{2}}{dx^{2}} f(x+tu)|_{t=o} = \langle \nabla^{2} f(x) u, u \rangle.$ 
That is,  $\langle \nabla^{2} f(x) u, u \rangle$  is the 2nd order derivative of  $f(x)$ 
along  $u$ .
  
• We can similarly show that
$$\frac{2}{\partial t_{1}} \frac{2}{\partial t_{n}} f(x+t_{1}u+t_{1}v)|_{t=t_{1}=0} = \frac{2}{\partial t_{n}} \frac{2}{\partial t_{1}} f(x+t_{1}u+t_{2}v)|_{t_{1}=t_{1}=0} = \frac{2}{\partial t_{n}} \frac{2}{\partial t_{n}} f(x+t_{1}u+t_{2}v)|_{t_{1}=t_{1}=0} = \frac{2}{\partial t_{n}} \frac{2}{\partial t_{n}} f(x+t_{1}v+t_{2}v)|_{t_{1}=t_{1}=0} = \frac{2}{\partial t_{n}} \frac{2}{\partial t_{n}} f(x+t_{1}v+t_{2}v)|_{t_{1}=t_{1}=t_{1}=t_{2}=0} = \frac{2}{\partial t_{n}} \frac{2}{\partial t_{n}} f(x+t_{1}v)|_{t_{1}=$$