

Final Review - ch4

Sunday, 6 December 2020 3:08 PM

Chapter 4: Linear functions and differentiation

Section 4.1: Linear Functions

Let $f: V \rightarrow \mathbb{R}$ be a function on a vector space V .
f is a linear function if
$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y) \quad \forall \alpha, \beta \in \mathbb{R} \text{ and } x, y \in V.$$

Example for linear function:
Mean function, $a^T x$, $F(f) = \int_a^b f(x) dx$, $f: L^1(a, b) \rightarrow \mathbb{R}$,
 $f(x) = \langle x, z \rangle$, $z \in V$

Example for non-linear function:
 $f(x) = \max\{x_1, x_2, \dots, x_n\}$, $g(x) = \|x\|$ (Norm)

Properties of linear functions

Properties of linear functions

- Homogeneity: $f(\alpha x) = \alpha f(x) \quad \forall \alpha \in \mathbb{R} \text{ and } x \in V.$
(Because $f(\alpha x) = f(\alpha x + 0x) = \alpha f(x) + 0 f(x) = \alpha f(x)$)
It implies $f(0) = 0$, because $f(0) = f(0 \cdot x) = 0 \cdot f(x) = 0 \quad \forall x \in V.$
- Additivity: $f(x+y) = f(x) + f(y) \quad \forall x, y \in V.$
- $f(\alpha_1 x_1 + \dots + \alpha_k x_k) = \alpha_1 f(x_1) + \dots + \alpha_k f(x_k), \quad \forall \alpha_1, \dots, \alpha_k \in \mathbb{R}, x_1, \dots, x_k \in V.$

To see this, we note that

$$\begin{aligned} f(\alpha_1 x_1 + \dots + \alpha_k x_k) &= f(\alpha_1 x_1) + f(\alpha_2 x_2 + \dots + \alpha_k x_k) \\ &= \alpha_1 f(x_1) + \alpha_2 f(x_2) + f(\alpha_3 x_3 + \dots + \alpha_k x_k) \\ &\vdots \\ &= \alpha_1 f(x_1) + \dots + \alpha_k f(x_k). \end{aligned}$$

Riesz representation theorem - Preliminary:

a linear function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ if and only if $f(x) = \langle a, x \rangle$ for some unique $a \in \mathbb{R}^n$

Pf14: any linear function can be written as inner product

Pf15: the inner product is unique

Riesz representation theorem - extend to hilbert space:

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

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

Examples for Riesz Representation Theorem:

1. Mean(x) = $\langle a, x \rangle$ where $a = (1/n, 1/n, \dots, 1/n)^T$
2. $F(f) = \int_a^b f(t) dt = \langle f, h \rangle$, where $h(t) = 1 \quad \forall t \in [a, b]$

Examples for not a Riesz Representation Theorem:

1. Norm (since norm is non linear)

$F: C[a, b] \rightarrow \mathbb{R}$ defined by $F(f) = f(t_0)$ for a fixed $t_0 \in \mathbb{R}$.

2. F is linear in $C[a, b]$

Pf16: a fixed t_0 function cannot have $h(t)$ in $\langle f, h \rangle$

Hyperplanes

The set $\{x \in V \mid \langle a, x \rangle = b\}$, where $a \in V$ and $b \in \mathbb{R}$ are given,
is called a Hyperplane in V .

Pf17: $S_{a, 0} = \{x \in V \mid \langle a, x \rangle = 0\}$ is a hyperplane.

Pf18: $S_{a, b} = \{x \in V \mid \langle a, x \rangle = b\}$ is also a hyperplane.

Projection onto hyperplanes

The vector on S that is the closest to y
is called the projection of y on S , denoted by $P_S y$,
i.e., $P_S y = \arg \min_{x \in S} \|x - y\|$.

Theorem: z is a solution of $\min_{x \in S} \|x - y\|$ if and only if
 $z \in S$ and $\langle z - y, z - x \rangle = 0$ if $x \in S$.



$(y - P_S y)$ vector should be orthogonal to x

Pf 19: proof of the above theorem

Theorem: The solution of $\min_{x \in S} \|x - y\|^2$ exists and unique, which is given by $y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$

Pf 20: $P_S y$ should be the above formula

- In summary, the projection $P_S y$ of $y \in V$ onto the hyperplane $S = \{x \in V \mid \langle a, x \rangle = b\}$ exists and is unique. Furthermore,

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

$$\langle P_S y - y, P_S y - x \rangle = 0.$$

Affine functions:

Affine functions

A linear function plus a constant is called an affine function.
That is, a function $f: V \rightarrow \mathbb{R}$ is affine if

$$f(x) = g(x) + b,$$
where $g: V \rightarrow \mathbb{R}$ is linear and $b \in \mathbb{R}$ is a constant.

Properties of affine functions:

Properties:

- If $f: V \rightarrow \mathbb{R}$ is affine, then linearity

Properties:

- If $f: V \rightarrow \mathbb{R}$ is affine, then linearity

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y) \quad \forall x, y \in V \text{ and } \alpha, \beta \in \mathbb{R} \text{ s.t. } \underline{\alpha + \beta = 1}.$$
To see this, $\begin{array}{l} \text{linear} \\ f(\alpha x + \beta y) = g(\alpha x + \beta y) + b = \alpha g(x) + \beta g(y) + (\alpha + \beta)b \\ = \alpha(g(x) + b) + \beta(g(y) + b) = \alpha f(x) + \beta f(y). \end{array}$

- If $f: V \rightarrow \mathbb{R}$, where V is a Hilbert space, then
 f must be in the form of Can be written as inner product

$$f(x) = \langle a, x \rangle + b, \text{ where } a \in V \text{ and } b \in \mathbb{R}.$$

Section 4.2: Case Studies: Regression and classification

Section 4.2.1: Linear Regression:

- Mathematically, we need to find a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f(x_i) \approx y_i, \quad i=1, 2, \dots, N$$

Linear model:

We search f in the class of all affine functions,
i.e., $f(x) = \langle a, x \rangle + b$ for some $a \in \mathbb{R}^n$, $b \in \mathbb{R}$.

Therefore, we find $a \in \mathbb{R}^n$, $b \in \mathbb{R}$ by solving

$$\min_{\substack{a \in \mathbb{R}^n \\ b \in \mathbb{R}}} \sum_{i=1}^N (\langle a, x_i \rangle + b - y_i)^2$$

再轉做下面的形式:

Write $X = \begin{bmatrix} x_1^T & | \\ x_2^T & | \\ \vdots & \vdots \\ x_N^T & | \end{bmatrix} \in \mathbb{R}^{N \times (n+1)}$, $\beta = \begin{bmatrix} a \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$
and $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$.

Then LS problem becomes

$$\min_{\beta \in \mathbb{R}^{n+1}} \|X\beta - y\|_2^2.$$

Since we have N linear equations to fit and $n+1$ unknowns,
 $N \geq n+1$ is required to have a unique solution.

Regularization:

好多時候我們地都不會有足夠的data, 即是 $N < n+1$

所以我們不會在 \mathbb{R}^{n+1} 裏面找 beta, 而是在一個subset裏面找

而這個subset的選定就有下列不同選擇:

1. Ridge regression

We choose $S = \{\beta = [\begin{matrix} a \\ b \end{matrix}] \in \mathbb{R}^{n+1} \mid \|a\|_2 \leq C\}$ for some $C > 0$.

Then we solve $\min_{\beta \in S} \|X\beta - y\|_2^2$

$\Downarrow \leftarrow$ by convex optimization theory
(Lagrangian)

$$\min_{\beta = [\begin{matrix} a \\ b \end{matrix}] \in \mathbb{R}^{n+1}} \frac{1}{2} \|X\beta - y\|_2^2 + \lambda \|a\|_2^2,$$

where $\lambda > 0$ depends on C and others.

Here $\|a\|_2^2$ is the regularization term.

and $\frac{1}{2} \|X\beta - y\|_2^2$ is the data-fitting term.

2. LASSO regression except above $\|a\|_2$ is replaced by $\|a\|_1$

Small $\|a\|_1$ tends to give a sparse vector $a \in \mathbb{R}^n$

(i.e., many entries of a are zeros)

Consequently, LASSO gives a solution $\beta = [\begin{matrix} a \\ b \end{matrix}]$ such that

$X\beta \approx y$ and a is sparse

Thus, given $X \in \mathbb{R}^{n \times p}$, the prediction is given by

$$\langle X, a \rangle + b = \sum_{i=1}^n a_i x_i + b \underset{\text{Let } I = \{i : a_i \neq 0\}}{\approx} \sum_{i \in I} a_i x_i + b$$

Since I is a small set, the prediction depends only on a small portion entries of X .

This is preferred because the prediction is interpretable.

Section 4.2.2 Kernel Regression

We solve

$$\min_{a \in H} \frac{1}{2} \sum_{i=1}^N (\langle a, \phi(x_i) \rangle - y_i)^2 + \lambda \|a\|_H^2$$

Representer Theorem:

The solution must be in the form of $a = \sum_{i=1}^N c_i \phi(x_i)$ for some $c = \begin{bmatrix} c_1 \\ \vdots \\ c_N \end{bmatrix} \in \mathbb{R}^N$

Pf21: proof of representer theorem

After some operation, we have the following result

Thus, the solution of the original minimization is

$$a = \sum_{i=1}^N c_i \phi(x_i)$$

where $c \in \mathbb{R}^N$ is a solution of $\min_{c \in \mathbb{R}^N} F_1(c)$.

$$F_1(c) = \frac{1}{2} \|Kc - y\|_2^2 + \lambda c^T K c$$

上述Function 等於下面function

$$\min_{c \in \mathbb{R}^N} \sum_{i=1}^N \left(\sum_{j=1}^N c_j K(x_i, x_j) - y_i \right)^2 + \lambda \sum_{i=1}^N \sum_{j=1}^N c_i c_j K(x_i, x_j)$$

等於先solve c 再solve a , 但係minimize兩個問題都係等價的

Then the predicted output y for input $x \in \mathbb{R}^n$ is

$$\forall x \in \mathbb{R}^n, f(\phi(x)) = y = \langle a, \phi(x) \rangle = \langle \sum_{i=1}^N c_i \phi(x_i), \phi(x) \rangle = \sum_{i=1}^N c_i K(x_i, x)$$

Summary: 整個kernel ridge regression algorithm等於

Kerel ridge regression.

Input: $(x_i, y_i), i=1, \dots, N$,
 $x_i \in \mathbb{R}^n, y_i \in \mathbb{R}$

Algorithm:

① Choose a kernel function $K(\cdot, \cdot)$,
e.g., $K(x, z) = e^{-\frac{\|x-z\|_2^2}{\sigma^2}}$, where $\sigma > 0$

② Solve $c \in \mathbb{R}^N$ by

$$\min_{c \in \mathbb{R}^N} \sum_{i=1}^N \left(\sum_{j=1}^N c_j K(x_i, x_j) - y_i \right)^2 + \lambda \sum_{i=1}^N \sum_{j=1}^N c_i c_j K(x_i, x_j)$$

Prediction: Given a new $x \in \mathbb{R}^n$, the predicted response is

$$\sum_{i=1}^N c_i K(x_i, x)$$

Pf22: prove how to change objective function to inner product only in kernel regression problem

Section 4.2.3 Classification

Problem Definition

- Classification: Giving training data

$(X_1, Y_1), (X_2, Y_2), \dots, (X_N, Y_N)$, $X_i \in \mathbb{R}^n$, $Y_i \in \{-1, +1\}$, $i=1, \dots, N$,
find a classifier (a function) f such that

$$y_i = \begin{cases} 1, & \text{if } f(X_i) \geq 1 \\ -1, & \text{if } f(X_i) \leq -1 \end{cases}$$

We use hyperplanes to separate the points

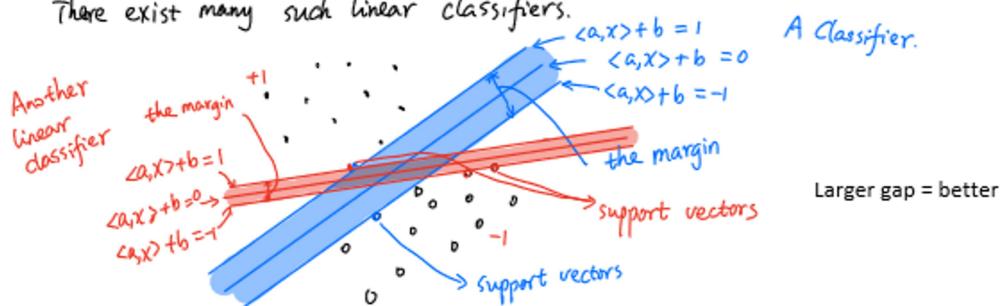
$$f(X) = \langle a, X \rangle + b, \text{ where } a \in \mathbb{R}^n, b \in \mathbb{R}.$$

The weights $a \in \mathbb{R}^n$ and $b \in \mathbb{R}$ are normalized such that

$$\langle a, X_i \rangle + b_i = \begin{cases} \geq 1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases}$$

- Support Vector Machine (SVM)

There exist many such linear classifiers.

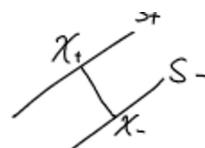


SVM-1:

Let us calculate the margin in terms of a and b .

The margin is the distance between the two hyperplanes

$$S_+ = \{x | \langle a, x \rangle + b = 1\} \text{ and } S_- = \{x | \langle a, x \rangle + b = -1\}$$



Let $X_+ \in S_+$ and $X_- \in S_-$ such that

$$\|X_+ - X_-\|_2 = \text{dist}(S_+, S_-).$$

個margin嘅寬度就係呢個

經過推論之後

$$\boxed{\begin{array}{ll} \min_{a \in \mathbb{R}^n, b \in \mathbb{R}} & \frac{1}{2} \|a\|_2^2 \\ \text{s.t.} & y_i (\langle a, X_i \rangle + b) \geq 1 \end{array}} \quad (\text{SVM-1})$$

Pf22: SVM的推論過程

Soft SVM:

We consider a "soft" version of (SVM-1) as follows.

We minimize the error to the separation

$$\sum_{i=1}^N h((y_i \langle a, x_i \rangle + b) - 1),$$

where the function $h(t) = \begin{cases} 0 & \text{if } t \geq 0 \\ |t| & \text{if } t < 0 \end{cases}$

SVM-2(Soft margin)

we relax the constraints to

$$\min_{\substack{a \in \mathbb{R}^n \\ b \in \mathbb{R} \\ \xi \in \mathbb{R}^N}} \frac{\lambda}{2} \|a\|_2^2 + \sum_{i=1}^N \xi_i \quad \text{s.t. } y_i(\langle a, x_i \rangle + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i=1, \dots, N$$

$$\Updownarrow \text{eliminate } \xi_i : \begin{cases} y_i(\langle a, x_i \rangle + b) - 1 \geq -\xi_i \\ \xi_i \geq 0 \end{cases}$$

$$\min_{\substack{a \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{\lambda}{2} \|a\|_2^2 + \sum_{i=1}^N h((y_i \langle a, x_i \rangle + b) - 1), \quad (\text{SVM-2})$$

where $h(t) = \begin{cases} 0 & \text{if } t \geq 0 \\ |t| & \text{if } t < 0 \end{cases}$

Kernel SVM Problem

Therefore, (SVM-2) becomes

$$\min_{a \in H} h(y_i \langle a, \phi(x_i) \rangle - 1) + \frac{\lambda}{2} \|a\|_H^2 \quad (K-SVM)$$

經過一陣推論後, (主要都係將 a 拆開 $a_s + \sum\{c_i \phi(x_i)\}$) Pf25: Kernel SVM Algo 推論

Thus, (K-SVM) becomes

$$\min_{C \in \mathbb{R}^N} h\left(y_i \left(\sum_{j=1}^N K(x_i, x_j) c_j\right) - 1\right) + \frac{\lambda}{2} C^T K C,$$

where $K = [K(x_i, x_j)]_{i=1, j=1}^{N, N} \in \mathbb{R}^{N \times N}$.

If we define $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ for some $K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$.

The prediction of the input x is given by

$$\text{sgn}\left(\sum_{j=1}^N K(x, x_j) c_j\right) \Rightarrow \text{就係用呢個黎 predict, 瞒佢大過 1 定細過 1}$$

Kernel SVM Algo Summary:

Kernel SVM:

① Choose a kernel function $K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$
 (e.g., $K(x, z) = e^{-\frac{\|x-z\|^2}{\sigma^2}}$)

② Solve C from

$$\min_{C \in \mathbb{R}^N} \frac{\lambda}{2} C^T K C + \sum_{i=1}^N h(y_i \cdot (Kc)_i - 1)$$

③ Then the classify function

Sections 4.3 Linear Approximation and differentiation

Definition of differentiation

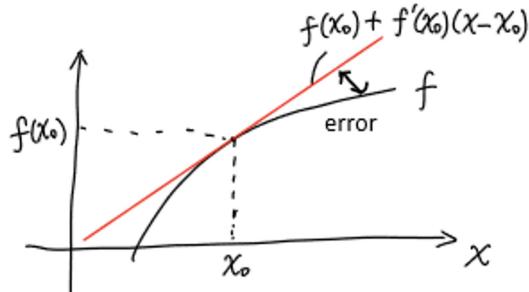
Recall that for a function $f: \mathbb{R} \rightarrow \mathbb{R}$, the derivative at x_0 is

$$f'(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0},$$

which is the same as

$$\lim_{x \rightarrow x_0} \left| \frac{f(x) - f(x_0) - f'(x_0)(x - x_0)}{x - x_0} \right| = 0.$$

Notice that $f(x_0) + f'(x_0)(x - x_0)$ is an affine function in \mathbb{R} that passes through $(x_0, f(x_0))$.



Differentiation 定義的轉換:

用 $f(x_0) + f'(x_0)(x - x_0)$ 去 approximate $f(x)$, 同埋從 differentiation 個定義來看, 可以得知

$$\lim_{x \rightarrow x_0} \left| \frac{\text{error}}{x - x_0} \right| \rightarrow 0 = \text{little } o \text{ of } (|x - x_0|)$$

Using this differentiation to extend to hilbert space, we can define Frechet differentiation

Definition: Let V be a Hilbert space. Let $f: V \rightarrow \mathbb{R}$. Then f is said Frechet differentiable if there exists a $v \in V$ such that

$$\lim_{x \rightarrow x^{(0)}} \frac{|f(x) - f(x^{(0)}) - \langle v, x - x^{(0)} \rangle|}{\|x - x^{(0)}\|} = 0. \quad \begin{pmatrix} \text{Note that} \\ x \rightarrow x^{(0)} \text{ is the} \\ \text{same as } \|x - x^{(0)}\| \rightarrow 0 \end{pmatrix}$$

If f is differentiable at $x^{(0)}$, v is called the gradient of f at $x^{(0)}$, denoted by $\nabla f(x^{(0)})$.

Example of how to calculate gradient using definition

Example 1: $f(x) = \|x\|^2$, where $\|x\|$ is the norm on V .

At any $x^{(0)} \in V$,

$$\begin{aligned}\|x\|^2 &= \|(x-x^{(0)})+x^{(0)}\|^2 = \langle (x-x^{(0)})+x^{(0)}, (x-x^{(0)})+x^{(0)} \rangle \\ &= \|x-x^{(0)}\|^2 + \|x^{(0)}\|^2 + 2\langle x^{(0)}, x-x^{(0)} \rangle\end{aligned}$$

Therefore, $\|x\|^2 - \underbrace{\left(\|x^{(0)}\|^2 + 2\langle x^{(0)}, x-x^{(0)} \rangle\right)}_{\text{affine approximation}} = \|x-x^{(0)}\|^2$

$$\text{So } \lim_{x \rightarrow x^{(0)}} \frac{\|x\|^2 - \|x^{(0)}\|^2 - 2\langle x^{(0)}, x-x^{(0)} \rangle}{\|x-x^{(0)}\|} = \lim_{x \rightarrow x^{(0)}} \frac{\|x-x^{(0)}\|^2}{\|x-x^{(0)}\|} = 0.$$

Thus, $\nabla f(x^{(0)}) = 2x^{(0)}$

目標就是要搬到 $\langle \text{grad}(x), x-x^{(0)} \rangle$ 這個形式

Properties: For frechet

① Frechet differentiation is linear, 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 $g \circ f: V \rightarrow \mathbb{R}$ and

$$\nabla(g \circ f)(x) = g'(f(x)) \cdot \nabla f(x)$$

if f and g are differentiable at x and $f(x)$ respectively.

Taylor expression

Taylor's expansion

From the definition, we see that

$$f(x) \approx f(x^{(0)}) + \langle \nabla f(x^{(0)}), x-x^{(0)} \rangle$$

Or, more precisely,

$$f(x) = f(x^{(0)}) + \langle \nabla f(x^{(0)}), x-x^{(0)} \rangle + o(\|x-x^{(0)}\|)$$

$$f(x) = f(x^{(0)}) + \langle \nabla f(x^{(0)}), x - x^{(0)} \rangle + o(\|x - x^{(0)}\|)$$

This is a generalization of Taylor's expansion.

In particular, if $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(x) \approx f(x^{(0)}) + \sum_{i=1}^n \frac{\partial f}{\partial x_i}(x^{(0)}) (x_i - x_i^{(0)})$$

Differentiation on normed vector spaces

Since no Riesz representation on normed spaces, only use linear function L in differentiation

Differentiation on normed vector spaces

Let V be a normed vector space.

Let $f: V \rightarrow \mathbb{R}$ be a function. Let $x^{(0)} \in V$.

To define differentiation, we still use an affine function approximation, and the affine function passes thru $(x^{(0)}, f(x^{(0)}))$.

Thus, we use $f(x^{(0)}) + L(x - x^{(0)})$, where $L: V \rightarrow \mathbb{R}$ is a linear function, to approximate $f(x)$.

However, because there is no Riesz representation on normed spaces, we keep the linear function L in the definition of differentiation.

Definition: f is differentiable at $x^{(0)} \in V$, if:

\exists a linear function $L: V \rightarrow \mathbb{R}$ such that

$$\lim_{\|x - x^{(0)}\| \rightarrow 0} \frac{|f(x) - (f(x^{(0)}) + L(x - x^{(0)}))|}{\|x - x^{(0)}\|} = 0.$$

The linear function L is called the differentiation of f at $x^{(0)}$.