

作业 2 证明题部分

1. [10 分] 试证：对于一组实数向量数据点的集合 $S = \{x_1, x_2, \dots, x_n\}$ ，若定义质心 c ，使得平方欧几里得距离和 $\sum_{i=1}^n \|x_i - c\|^2$ 最小，则质心 c 就是所有点的均值 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ 。

(a) 若用欧氏距离，则质心为所有点的均值。

Proof: 设 x_i 坐标为 $(x_{i1}, x_{i2}, \dots, x_{im})$, $i=1 \dots n$.

c 坐标为 (c_1, \dots, c_m) .

$$D(c) = \sum_{i=1}^n \|x_i - c\|^2 = \sum_{j=1}^m \sum_{i=1}^n (x_{ij} - c_j)^2.$$

$$D'(c) = -2 \sum_{j=1}^m \sum_{i=1}^n (x_{ij} - c_j).$$

$$D''(c) = 2mn > 0.$$

\Rightarrow 当 $D'(c)=0$ 时 $D(c)$ 有 min.

$$\text{令 } D'(c)=0 \text{ 得 } c_j = \frac{1}{n} \sum_{i=1}^n x_{ij} \quad j=1, 2, \dots, m$$

即当 $c = \bar{x}$ 时平方欧几里得距离最小。

\Rightarrow 质心为所有点的均值 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

1. [10 分] 试证：对于一组实数向量数据点的集合 $S = \{x_1, x_2, \dots, x_n\}$ ，若定义质心 c ，使得曼哈顿距离和 $\sum_{i=1}^n |x_i - c|$ 最小，则质心 c 就是所有点的中位数。

(b) 带用曼哈顿距离，则质心为所有点中位数。

Proof: 设 x_i 坐标为 $(x_{i1}, x_{i2}, \dots, x_{im})$ $i=1 \dots n$.

c 坐标为 (c_1, c_2, \dots, c_m) .

$$D(c) = \sum_{i=1}^n |x_{i1} - c_1| \geq 0.$$

先看第一维的情况：

$$\text{令 } d_1(c) = \sum_{i=1}^n |x_{i1} - c_1|.$$

$\because |x_{i1} - c_1|$ 为凸函数 $\Rightarrow d_1(c)$ 也是凸函数

$\Rightarrow d_1(c)$ 的最小值存在，且取 min 时 $d'_1(c) = 0$.

由 $S = \{x_1, \dots, x_n\}$ 的无序性，不妨令 $x_1 \geq x_2 \geq \dots \geq x_n$.

再设 $x_{k1} \geq c_1 \geq x_{(k+1)1}$

$$\text{则 } d_1(c) = \sum_{i=1}^k (x_{i1} - c_1) - \sum_{i=k+1}^n (x_{i1} - c_1).$$

$$d'_1(c) = -k + (n-k) = 0.$$

$$c = \frac{n}{2}$$

当 n 为偶数时， $x_{\frac{n}{2}1} \geq c_1 \geq x_{(\frac{n}{2}+1)1} \Rightarrow c_1$ 为该区间内任一点。

当 n 为奇数时，只有当 $c_1 = x_{\frac{n+1}{2}}$ 时能使 $d'_1(c) = 0$. 此时 c_1 亦为中位数。

综合奇偶，只有当 c_1 为中位数时能使得 $d_1(c)$ 达到最小值，即 0.

同理，对于其他维度，也是 c_2, \dots, c_m 取中位数时让 $d_2(c), \dots, d_m(c)$ 为 0

此时 $D(c) = \sum_{j=1}^m d_j(c) = 0$ 取到最小值。

\Rightarrow 质心 c 为所有点的中位数。