

Introduction to Machine Learning, Fall 2015 - Exercise session VI

Rodion “rodde” Efremov, student ID 013593012

December 1, 2015

Problem 3 (3 points)

Given a set of N points y_1, \dots, y_N , with each $y_i \in \mathbb{R}$, show that

(a) the value y^* which minimizes the sum of *squared* errors, i.e.

$$y^* = \arg \min_{\hat{y}} \sum_{i=1}^N (y_i - \hat{y})^2$$

is given by the *mean* of the y_i , i.e. $y^* = \sum_i y_i / N$.

(b) the value y^* which minimizes the sum of *absolute* errors, i.e.

$$y^* = \arg \min_{\hat{y}} \sum_{i=1}^N |y_i - \hat{y}|$$

is given by the *median* of the y_i .

Solution to (a)

Let us first rewrite the sum:

$$f(x) = \sum_{i=1}^N (y_i - x)^2 = \sum_{i=1}^N (y_i^2 - 2y_i x + x^2).$$

We want to find the value of x which will minimize $f(x)$. Derivating $f(x)$ with respect to x :

$$\frac{d}{dx} f(x) = \sum_{i=1}^N (-2y_i + 2x),$$

since derivation is a linear operator. We will also need the second derivative of f :

$$\frac{d^2}{dx^2} f(x) = \sum_{i=1}^N 2 = 2N.$$

Since $N > 0$, we have that the second derivative is always positive. Also, a root of the first derivative of f is given by equation

$$\sum_{i=1}^N (2x - 2y_i) = 0,$$

which leads to

$$\begin{aligned}\sum_{i=1}^N 2x &= \sum_{i=1}^N 2y_i, \\ 2Nx &= \sum_{i=1}^N 2y_i, \\ Nx &= \sum_{i=1}^N y_i, \\ x &= \frac{1}{N} \sum_{i=1}^N y_i,\end{aligned}$$

which is a requested average of y_1, \dots, y_N . Note since the second derivative is always positive, f attains its **minimum** at x .

Solution to (b)

Suppose that the values y_1, \dots, y_N are **sorted** ($y_1 \leq y_2 \leq \dots \leq y_N$). Now consider the smallest and the largest values: y_1 and y_N , respectively. Denote by x the value that minimizes the sum of absolute errors. As long $x \in [y_1, y_N]$, the sum $|x - y_1| + |y_N - x|$ is minimized (if you take x outside of $[y_1, y_N]$, you will increase the value of $|x - y_1| + |y_N - x|$ beyond its minimal possible value). Now, in order to minimize the sum of absolute errors, we **must** keep x within the range $[y_1, y_N]$. Now consider the second smallest and the second largest values: y_2 and y_{N-1} , respectively. Now since $y_2 \geq y_1$ and $y_{N-1} \leq y_N$, keeping x within the range $[y_1, y_N]$ is not enough as x might be outside of the range $[y_2, y_{N-1}]$ which will increase the sum of absolute errors above its minimum. Using the same reasoning, we continue this until reaching a median (or two of them).

If N is even, the strongest bound for x is $[y_{N/2}, y_{N/2+1}]$. We can choose

$$x = \frac{y_{N/2} + y_{N/2+1}}{2},$$

which is, in fact, the median of $(y_i)_{i=1}^N$. However, if N is odd, we must satisfy $x \in [y_{(N+1)/2}, y_{(N+1)/2}]$, which is satisfied only if $x = y_{(N+1)/2}$, which is a median once again.