

Lecture Notes II.1 – Bias and variance in Kernel Regression

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

October, 2023

An elementary analysis

Bias, Variance and h for $x \in \mathbb{R}$

Kernel regression by Nadaraya-Watson

$$\hat{y}(x) = \frac{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right) y^i}{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right)} \quad (1)$$

$$\text{Let } w_i = \frac{b\left(\frac{\|x-x^i\|}{h}\right)}{\sum_{i'=1}^n b\left(\frac{\|x-x^{i'}\|}{h}\right)}.$$

Assumptions

A0 For simplicity, in this analysis we assume $x \in \mathbb{R}$.

A1 There is a true smooth¹ function $f(x)$ so that

$$y = f(x) + \varepsilon, \quad (2)$$

where ε is sampled independently for each x from a distribution P_ε , with $E_{P_\varepsilon}[\varepsilon] = 0$, $\text{Var}_{P_\varepsilon}(\varepsilon) = \sigma^2$.

A2 The kernel $b(z)$ is smooth, $\int_{\mathbb{R}} b(z) dz = 1$, $\int_{\mathbb{R}} zb(z) dz = 0$, and we denote $\sigma_b^2 = \int_{\mathbb{R}} z^2 b(z) dz$, $\gamma_b^2 = \int_{\mathbb{R}} b^2(z) dz$.

In this first analysis, we consider that the values x , $x^{1:N}$ are fixed; hence, the randomness is only in $\varepsilon^{1:N}$.

¹with continuous derivatives up to order 2

Expectation of $\hat{y}(x)$ – a simple analysis

Expanding f in Taylor series around x we obtain

$$f(x^i) = f(x) + f'(x)(x^i - x) + \frac{f''(x)}{2}(x^i - x)^2 + o((x^i - x)^2) \quad (3)$$

We also have

$$y^i = f(x^i) + \varepsilon^i. \quad (4)$$

We now write the expectation of $\hat{y}(x)$ from (1), replacing in it y^i and $f(x^i)$ as above. What we would like to happen is that this expectation equals $f(x)$. Let us see if this is the case.

$$E_{P_{\varepsilon}^n} [\hat{y}(x)] = E_{P_{\varepsilon}^n} \left[\sum_{i=1}^n w_i y^i \right] = E_{P_{\varepsilon}^n} \left[\sum_{i=1}^n w_i (f(x^i) + \varepsilon^i) \right] \quad (5)$$

$$= \sum_{i=1}^n w_i f(x) + \sum_{i=1}^n w_i f'(x)(x^i - x) + \sum_{i=1}^n w_i \frac{f''(x)}{2}(x^i - x)^2 + \underbrace{E_{P_{\varepsilon}^n} \left[\sum_{i=1}^n w_i \varepsilon^i \right]}_{=0} \quad (6)$$

$$= \underbrace{f(x) + f'(x) \sum_{i=1}^n w_i (x^i - x) + \frac{f''(x)}{2} \sum_{i=1}^n w_i (x^i - x)^2}_{\text{bias}} \quad (7)$$

In the above, the expressions in red depend of f , those in blue depend on x and $x^{1:n}$.

Qualitative analysis of the bias terms

The first order term $f'(x) \sum_{i=1}^n w_i (x^i - x)$ is responsible for **border effects**.

The second order term **smooths out** sharp peaks and valleys.

Bias, Variance and h for $x \in \mathbb{R}$

2

The **bias** of \hat{y} at x is defined as $E_{P_X^n} E_{P_\varepsilon^n} [\hat{y}(x) - f(x)]$.

$$E_{P_X^n} E_{P_\varepsilon^n} [\hat{y}(x) - f(x)] = h^2 \sigma_b^2 \left(\frac{f'(x) p_X'(x)}{p_X(x)} + \frac{f''(x)}{2} \right) + o(h^2) \quad (8)$$

The **variance** \hat{y} at x is defined as $\text{Var}_{P_X^n P_\varepsilon^n}(\hat{y}(x))$.

$$\text{Var}_{P_X^n P_\varepsilon^n}(\hat{y}(x)) = \frac{\gamma^2}{nh} \sigma^2 + o\left(\frac{1}{nh}\right). \quad (9)$$

The **MSE (Mean Squared Error)** is defined as $E_{P_X^n} E_{P_\varepsilon^n} [(\hat{y}(x) - f(x))^2]$, which equals

$$\text{MSE}(x) = \text{bias}^2 + \text{variance} = h^4 \sigma_b^4 \left(\frac{f'(x) p_X'(x)}{p_X(x)} + \frac{f''(x)}{2} \right) + \frac{\gamma_b^2}{nh} \sigma^2 + \dots \quad (10)$$

Optimal selection of h

If the MSE is integrated over \mathbb{R} we obtain the **MISE** $= \int_{\mathbb{R}} \text{MSE}(x) p_X(x) dx$.
The kernel width h can be chosen to minimize the MISE, for fixed f, p_X and b .
We set to 0 the partial derivative

$$\frac{\partial \text{MISE}}{\partial h} = h^3 \left(\text{[red box]} \right) - \frac{\text{[blue box]}}{nh^2} = 0. \quad (11)$$

It follows that $h^5 \propto \frac{1}{n}$, or

$$h \propto \frac{1}{n^{1/5}}. \quad (12)$$

In d dimensions, the optimal h depends on the sample size n as

$$h \propto \frac{1}{n^{1/(d+4)}}. \quad (13)$$