



Lecture 2 (tablet, finally !!)

Read: All of NP Statisfics 5.1-5.4 Old ll posted LI, LII posted

large enough ?

Lecture II - Nearest Neighbor and Kernel predictors

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

STAT/BIOST 527 Spring 2023 1 Nearest-Neighbor predictors



In elementary analysis of Kernel Regression



(4) Bias, Variance and **h** for $x \in \mathbb{R}$

The Nearest-Neighbor predictor

- 1-Nearest Neighbor The label of a point x is assigned as follows:
 - **(**) find the example x^i that is nearest to x in \mathcal{D} (in Euclidean distance)
 - assign x the label yⁱ, i.e.

 $\hat{y}(x) = y^i$

k-NN for discrete x, y E R. define distance d(X,X)

The Nearest-Neighbor predictor

- 1-Nearest Neighbor The label of a point x is assigned as follows:
 - **(**) find the example x^i that is nearest to x in \mathcal{D} (in Euclidean distance)
 - 2 assign x the label y^i , i.e.

 $\hat{y}(x) = y^i$

- K-Nearest Neighbor (with K = 3, 5 or larger)
 - find the K nearest neighbors of x in D: xⁱ1,...iK
 for classification f(x) = the most frequent label among the K neighbors (well suited for multiclass)
 - for regression $f(x) = \frac{1}{K} \sum_{i \text{ neighbor of } x} y^i$ = mean of neighbors' labels



The Nearest-Neighbor predictor

- 1-Nearest Neighbor The label of a point x is assigned as follows:
 find the example xⁱ that is nearest to x in D (in Euclidean distance)
 - 2 assign x the label y^i , i.e.

 $\hat{y}(x) = y^i$

- K-Nearest Neighbor (with K = 3, 5 or larger)
 - find the K nearest neighbors of x in \mathcal{D} : x^{i_1, \dots, i_K}
 - for classification f(x) = the most frequent label among the K neighbors (well suited for multiclass)
 - for regression $f(x) = \frac{1}{K} \sum_{i \text{ neighbor of } x} y^i$ = mean of neighbors' labels

- No parameters to estimate!
- No training!
- But all data must be stored (also called memory-based learning)

Kernel regression and classification

igation. values in data Like the K-nearest neighbor but with "smoothed" neighborhoods • The predictor prediction: $f(x) = \sum_{i=1}^{n} \beta_i b(x, \underline{x}^i) y^i$ weight depend on X¹, X⁽¹⁾

where β_i are coefficients

new x



WiyeRh $f(x) = \sum W_i(x) y_i$ = 1 for all X (=) y = weighted aug $\sum_{i=1}^{n} W_i(x) = 1$ $f(x') = W(x')^T Y$ $f(x^i) = W(x^{1:n}) \mathcal{Y}$ Linear in \mathcal{Y}

Kernel regression and classification

- Like the K-nearest neighbor but with "smoothed" neighborhoods
- The predictor

$$f(x) = \sum_{i=1}^{n} \beta_i b(x, x^i) y^i$$
 (1)

where β_i are coefficients

- Intuition: center a "bell-shaped" kernel function b on each data point, and obtain the prediction f(x) as a weighted sum of the values y^i , where the weights are $\beta_i b(x, x^i)$
- Requirements for a kernel function b(x, x')
 - In non-negativity
 - 2 symmetry in the arguments x, x'
 - Optional: radial symmetry, bounded support, smoothness
- A typical kernel function is the Gaussian kernel (or Radial Basis Function (RBF))

$$b(z) \propto e^{-z^2/2} \tag{2}$$

$$b_h(x, x') \propto e^{-\frac{||x-x'||^2}{2h^2}}$$
 with $h =$ the kernel width (3)

Regression example

A special case in wide use is the Nadaraya-Watson regressor

In this regressor, f(x) is always a convex combination of the y^{i} 's, and the weigths are proportional to $b_h(x, x^i)$.

The Nadaraya-Watson regressor is biased if the density of P_X varies around x.



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

(4)

An example: noisy data from a parabola

STAT/BIOST 527 Spring 2023



Local Linear Regression

To correct for the bias (to first order) one can estimate a regression line around x.

- **O** Given query point x
- **2** Compute kernel $b_h(x, x^i) = w_i$ for all i = 1, ..., N
- Solve weighted regression $\min_{\beta,\beta_0} \sum_{i=1}^d w_i (y^i \beta^T x^i \beta_0)^2$ to obtain β,β_0 (β,β_0 depend on x through w_i)
- $Calculate f(x) = \beta^T x + \beta_0$

Exercise Show that Nadaraya-Watson solves a local linear regression with fixed $\beta = 0$

Kernel binary classifiers

- Obtained from Nadaraya-Watson by setting y^i to ± 1 .
- Note that the classifier can be written as the difference of two non-negative functions

 $f(x) \propto \sum_{i \neq i=1}^{h} b\left(\frac{||x-x^i||}{h}\right) \sum_{i \neq i=1}^{h} b\left(\frac{||x-x^i||}{h}\right).$ (5)2 = 1 (x', y'), i= 1: Ny y'e1±13 classifier $f(x) \in \frac{1}{2} + \frac{1}{2} = \frac{1}{2} \hat{y}(x) = sgnf(x)$ $f(x) \in [-1, +1] \sim \text{confidence}$ OR BiZA

Kernel binary classifiers

- Obtained from Nadaraya-Watson by setting y^i to ± 1 .
- Note that the classifier can be written as the difference of two non-negative functions

$$f(x) \propto \sum_{i:y^i=1} b\left(\frac{||x-x^i||}{h}\right) - \sum_{i:y^i=-1} b\left(\frac{||x-x^i||}{h}\right).$$
(5)



Kernel regression by Nadaraya-Watson

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

(6)

.

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

Assumptions

A0 For simplicity, in this analysis we assume $x \in \mathbb{R}$. A1 There is a true smooth¹ function f(x) so that

$$y = f(x) + \varepsilon, \tag{7}$$

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

A2 The kernel b(z) is smooth, $\int_{\mathbb{R}} b(z)dz = 1$, $\int_{\mathbb{R}} zb(z) = 0$, 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
subimator

Expanding f in Taylor series around x we obtain

$$f(x^{i}) = f(x) + f'(x)(x^{i} - x) + \frac{f''(x)}{2}(x^{i} - x)^{2} + o((x^{i} - x)^{2})$$
(8)

We also have

Wanter (x) =f(x) UNBIASED

$$y^{i} = f(x^{i}) + \varepsilon^{i}.$$
 (9)

We now write the expectation of $\hat{y}(x)$ from (6), 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.

$$E_{P_{\varepsilon}^{n}}[\hat{y}(x)] = E_{P_{\varepsilon}^{n}}\left[\sum_{i=1}^{n} w_{i}y^{i}\right] = E_{P_{\varepsilon}^{n}}\left[\sum_{i=1}^{n} w_{i}\left(f(x^{i}) + \varepsilon^{i}\right)\right]$$
(10)
$$= \sum_{i=1}^{n} w_{i}f(x) + \sum_{i=1}^{n} w_{i}f'(x)(x^{i} - x) + \sum_{i=1}^{n} w_{i}\frac{f''(x)}{2}(x^{i} - x)^{2} + \underbrace{E_{P_{\varepsilon}^{n}}\left[\sum_{i=1}^{n} w_{i}\varepsilon^{i}\right]}_{=0]}$$
(11)
$$= f(x) + f'(x)\sum_{i=1}^{n} w_{i}(x^{i} - x) + \frac{f''(x)}{2}\sum_{i=1}^{n} w_{i}(x^{i} - x)^{2}$$
(12)
bias

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

STAT/BIOST 527 Spring 2023

Qualitative analysis of the bias terms

The first order term f'(x) ∑_{i=1}ⁿ w_i(xⁱ − x) is responsible for border effects.
 The second order term smooths out sharp peaks and valleys.



Qualitative analysis of the bias terms

The first order term f'(x) ∑_{i=1}ⁿ w_i(xⁱ − x) is responsible for border effects.
 The second order term smooths out sharp peaks and valleys.



Qualitative analysis of the bias terms

The first order term f'(x) ∑_{i=1}ⁿ w_i(xⁱ − x) is responsible for border effects.
 The second order term smooths out sharp peaks and valleys.



and Kernel

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

2

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

$$\mathsf{E}_{\mathsf{P}_{X}^{n}}\mathsf{E}_{\mathsf{P}_{\varepsilon}^{n}}[\hat{y}(x) - f(x)] = h^{2}\sigma_{b}^{2}\left(\frac{f'(x)p'_{X}(x)}{p_{X}(x)} + \frac{f''(x)}{2}\right) + o(h^{2})$$
(13)

The variance \hat{y} at x is defined as $Var_{P_{x}^{n}}P_{\varepsilon}^{n}(\hat{y}(x))$.

$$Var_{P_{\chi}^{n}}P_{\varepsilon}^{n}(\hat{y}(\chi)) = \frac{\gamma^{2}}{nh}\sigma^{2} + o\left(\frac{1}{nh}\right).$$
(14)

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

$$MSE(x) = \text{bias}^{2} + \text{variance} = h^{4}\sigma_{b}^{4} \left(\frac{f'(x)p'_{X}(x)}{p_{X}(x)} + \frac{f''(x)}{2}\right) + \frac{\gamma_{b}^{2}}{nh}\sigma^{2} + \dots$$
(15)

Optimal selection of h

If the MSE is integrated over \mathbb{R} we obtain the MISE= $\int_{\mathbb{R}} MSE(x)p_X(x)dx$. 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

$$\frac{\partial MISE}{\partial h} = h^3 \left(\underbrace{}_{nh^2} \right) - \underbrace{\left(\underbrace{}_{nh^2} \right)}_{nh^2} = 0.$$
(16)
It follows that $h^5 \propto \frac{1}{n}$, or
 $h \propto \frac{1}{n^{1/5}}.$ (17)

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

$$h \propto \frac{1}{n^{1/(d+4)}}.$$
 (18)

Marina Meila: LII k-NN and Kernel

The MSE with optimal *h* decreases as
$$\sim \frac{1}{n}n^{1/(d+4)}$$

Compare this with the MSE of the mean of a distribution, which decreases $\sim \frac{1}{n}$