

A short notes on Nonparametric regression

Tong Zhou*

November 1, 2020

This notes is mainly based on the following materials:

- All of Nonparametric Statistics, by *Larry Wasserman*
- A Primer on Regression Splines, by *Jeffrey S. Racine*
- Lecture 3: Regression: Nonparametric Approaches, STAT 535 SML Autumn 2019, by *Yen-chi Chen*
- Lecture notes: Nonparametric Regression and Classification, Statistical Machine Learning, Spring 2017, by *Ryan Tibshirani and Larry Wasserman*

Given pairs of data $(x_1, Y_1), \dots, (x_n, Y_n)$, in which Y is called **response variable** and x is called **covariate**, our goal is to relate them by:

$$Y_i = r(x_i) + \varepsilon_i, \quad \mathbf{E}[\varepsilon_i] = 0, i = 1, \dots, n.$$

where $r()$ is called the **regression function**. In the machine learning community, the goal is to “learn” r under weak assumptions. The estimator of r is denoted by $\hat{r}_n(x)$, which is often referred to as a **smoother**.

Note that at this point, the interpretations regarding this model vary depending on different assumptions imposed on the data. Let’s discuss them:

1. Assume the random pair $(X, Y) \in \mathcal{P}$ and $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}, i = 1, \dots, n$ are i.i.d. . The regression function is **defined** as $r(x) = \mathbf{E}[Y|X = x]$. Then ε_i can be perceived as the error and equal to $\varepsilon_i = y_i - r(x_i)$. Under this setup, it is natural to get $\mathbf{E}[\varepsilon_i] = 0$ by the *Law of Iterated Expectations* and ε_i are i.i.d., for any i . This interpretation agrees with the one in Bruce’s Hansen’s Econometric textbooks.

*Johns Hopkins University, tzhou11@jhu.edu

2. Assume $\varepsilon_i \perp\!\!\!\perp X_i$. This independence assumption is not free. Bear in mind that throughout the discussion, we maintain the strong assumption. Why this assumption is strong? Think about a multiple regression example, if a relevant regressor is omitted which is also correlated with some existing regressor, then this assumption fails to hold. Specifically, if $y_i = \beta X_i + \alpha Z_i + \varepsilon_i$ and $\varepsilon_i \perp\!\!\!\perp (X_i, Z_i)$. If we adhere to the regression $y_i = \beta X_i + \tilde{\varepsilon}_i$ and X_i and Z_i happen to not be independent, then there is no hope $X_i \perp\!\!\!\perp \tilde{\varepsilon}_i$.
3. **Fixed x world.** $E[\varepsilon_i] = 0$ and $E[\varepsilon_i^2] = \sigma^2$. In this scenario, the *Fixed* inputs assumption is **equivalent to** the *Random* inputs with the independence assumption maintained.

We have mentioned r is called a smoother. Before plunging into the formal treatment of nonparametric regression, we first define the notion of **linear smoother**.

Definition 1. An estimator \hat{r}_n of r is a **linear smoother** if, for each x , there exists a vector $\ell(x) = (\ell_1(x), \dots, \ell_n(x))^T$ such that

$$\hat{r}_n(x) = \sum_{i=1}^n \ell_i(x) Y_i.$$

Some related quantities are:

- **fitted values:** $\mathbf{r} = (\hat{r}_n(x_1), \dots, \hat{r}_n(x_n))^T$
- **Smoothing matrix / hat matrix:** $L = (\ell(x_1), \dots, \ell(x_n))^T$ and $L_{ij} = \ell_j(x_i)$.
- **effective kernel** for $r(x_i)$: $\ell(x_i)^T$ the i -th row of L with $\sum_{i=1}^n \ell_i(x) = 1$.
- **effective degrees of freedom:** $\nu = \text{tr}(L)$.
- $\mathbf{r} = L\mathbf{Y}$.

NOTE

For the regression model:

$$Y_i = r(x_i) + \varepsilon_i$$

The estimator \hat{r}_n of r is to minimize the sum of squares:

$$\sum_{i=1}^n (Y_i - \hat{r}_n(x_i))^2.$$

There are two extreme cases:

1. Minimizing over *ALL linear functions* \implies **least square estimators.**
2. Minimizing over *ALL functions* \implies **a function that interpolates the data.**

There are several solutions:

1. Local regression: use a locally weighted sums of squares.

$$\sum_{i=1}^n w_i(x) (Y_i - a)^2.$$

2. Penalized sums of squares :

$$M(\lambda) = \sum_i (Y_i - \hat{r}_n(x_i))^2 + \lambda J(r)$$

Example 1 (Local Regression). Suppose $x_i \in \mathbb{R}$ is scalar, and consider the regression model $Y_i = r(x_i) + \varepsilon_i$. Idea of estimating $r(x)$ is by taking a weighted average of Y_i , giving higher weight to those points near x .

– **Nadaraya-Watson kernel estimator:**

$$\hat{r}_n(x) = \sum_{i=1}^n \ell_i(x) Y_i.$$

– K is a kernel and

$$\ell_i(x) = \frac{K\left(\frac{x-x_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)}$$

- *Example 5.26 (boxcar kernel:)* Suppose that $x_i \in [a, b], i = 1, \dots, n$. Divide (a, b) into m equally spaced bins denoted by B_1, B_2, \dots, B_m . Define $\hat{r}_n(x)$ by

$$\hat{r}_n(x) = \frac{1}{k_j} \sum_{i: x_i \in B_j} Y_i, \text{ for } x \in B_j$$

where k_j is the number of points in B_j . So for $x \in B_j$, $\ell_i(x) = \frac{1}{k_j}$ if $x_i \in B_j$ and $\ell_i = 0$ otherwise.

Example 2 (Penalized Regression). Goal: to minimize the *penalized sums of squares*:

$$M(\lambda) = \sum_i (Y_i - \hat{r}_n(x_i))^2 + \lambda J(r)$$

where $J(r)$ is some *roughness penalty*, λ controls the amount of smoothing.

1. $\lambda \rightarrow 0$: the interpolating function.
2. $\lambda \rightarrow \infty$: the least squares line.
3. $0 < \lambda < \infty$: use *splines*, a special piecewise polynomial.

A special case on the penalty term is such that

$$J(r) = \int |r''(x)|^2 dx \tag{1}$$

Before diving into details, we first look at an important theorem:

Theorem 1. The function $\hat{r}_n(x)$ that minimizes $M(\lambda)$ with penalty defined in 1 is a natural cubic spline with knots at the data points. The estimator \hat{r}_n is called a **smoothing spline**.

Now let us get back to some notions underlined in this theorem:

- Let $\xi_1 < \xi_2 < \dots < \xi_p$ be set of ordered points-called **knots**.
- **cubic spline**: is a continuous function r such that (i) r is a cubic piecewise polynomial over each bins $(\xi_1, \xi_2), \dots$ and (ii) r has continuous first and second derivatives at the knots.

However, Theorem 1 only asserts that the minimizer of the penalty regression problems exists and belongs to the class of cubic splines. It does not provide guidance on what the cubic spline looks like. The next theorem fills a gap:

Theorem 2. Let $\xi_1 < \xi_2 < \dots < \xi_p$ be knots contained in an interval (a, b) . Define $h_1(x) = 1, h_2(x) = x, h_3(x) = x^2, h_4(x) = x^3, h_{j+4}(x) = (x - \xi_j)_+^3$ for $j = 1, \dots, p$. The functions $\{h_1, \dots, h_{k+4}\}$ form a basis for the set of cubic splines at these knots, called the **truncated power basis**. Thus, any cubic spline $r(x)$ with these knots can be written as

$$r(x) = \sum_{j=1}^{k+4} \beta_j h_j(x).$$

Remark 1. – The space of the 3-order splines with knots ξ_1, \dots, ξ_p has dimension $p + 4$.

- The basis of functions is called **truncated power basis**.
- More general, the k -order spline r is a piecewise polynomial function of degree k that is continuous and has continuous derivatives of orders $1, \dots, k - 1$, at its knot points. Specifically, its corresponding *truncated power basis* is $g_1, g_2, \dots, g_{p+k+1}$ such that

$$g_1(x) = 1, g_2(x) = x, \dots, g_{k+1}(x) = x^k, \\ g_{k+1+j}(x) = (x - \xi_j)_+^k, j = 1, \dots, p.$$

and its dimension is $p + k + 1$.

To gain more computational edge, another more popular spline is **B-spline**, where its basis functions are defined recursively.

Let $\xi_0 = a$ and $\xi_{p+1} = b$. Define the appended $2M$ knots by:

$$\tau_1 \leq \tau_2 \leq \tau_3 \leq \dots \leq \tau_M \leq \xi_0,$$

and

$$\xi_{p+1} \leq \tau_{k+M+1} \leq \dots \leq \tau_{p+2M}$$

This appending is needed due to the recursive nature of the B-spline. Now for the j -th order of spline, the basis functions can be defined as:

$$B_{i,1}(x) = \begin{cases} 1 & \text{if } \tau_i \leq x \leq \tau_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

for each i . Next, for $m \leq M$, we define

$$B_{i,m}(x) = \alpha_{i,m-1} B_{i,m-1}(x) + (1 - \alpha_{i+1,m}) B_{i+1,m}(x)$$

where

$$\alpha_{i,m-1}(x) = \begin{cases} \frac{x - \tau_i}{\tau_{i+m-1} - \tau_i}, & \text{if } \tau_{i+m-1} \neq \tau_i \\ 0, & \text{otherwise} \end{cases}$$

Now all relevant notations are squared away, the following theorem dictates the role of B-spline:

Theorem 3. The function $\{B_{i,4}, i = 1, \dots, k + 4\}$ are a basis for the set of cubic splines. They are called the **B-spline basis functions**.

Let's switch back to the original setting: $(x_1, Y_1), \dots, (x_n, Y_n)$. Then \hat{r}_n is a natural cubic spline:

$$\hat{r}_n(x) = \sum_{j=1}^N \hat{\beta}_j B_j(x)$$

where $N = n + 4$ and $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_N)^\top$.

The minimization problem then can be expressed as

MINIMIZATION OF PENALIZED REGRESSION ESTIMATOR

$$\text{minimize: } (Y - B\beta)^\top (Y - B\beta) + \lambda \beta^\top \Omega \beta$$

where $B_{ij} = B_j(X_i)$ and $\Omega_{jk} = \int B_j''(x) B_k''(x) dx$

The value of β that minimizes the problem is:

$$\hat{\beta} = (B^\top B + \lambda \Omega)^{-1} B^\top Y.$$

Curse of dimensionality

If f_0 is L -Lipschitz continuous, the local polynomial estimator's error rate is $n^{-2/(2+d)}$, i.e.

$$\mathbb{E} \left[\|\hat{f} - f_0\|_2^2 \right] \lesssim n^{-2/(2+d)}.$$

It exhibits a very discouraging result: the error rate depends heavily on the dimension d . To be more concrete, fix a tolerance rate $\varepsilon > 0$, how large is the sample size n that can ensure $n^{-2/(2+d)} \leq \varepsilon$? By simple rearrangement, we have

$$n \geq \varepsilon^{-(2+d)/2}.$$

That is, as d is increasing, n is required to exponentially increase to achieve an error bound of ε . Recall in the linear regression models, we have a more reasonable requirement for the sample size $n \geq \frac{d}{\varepsilon}$.

To make it more precise, consider a minimax problem, which asserts that we cannot hope to do better than the error rate $n^{-2/(2+d)}$ over all L -Lipschitz function in d dimensions, denoted by $H_d(1, L)$. That is

$$\inf_{\widehat{f}} \sup_{f_0 \in H_d(1, L)} \mathbf{E} \left[\|\widehat{f} - f_0\|_2^2 \right] \gtrsim n^{-2/(2+d)}$$

So to circumvent this curse, more assumptions about it are to be made. One such example is the additive model.

Exercises

a.

Assume that we observe i.i.d. samples $(x_i, y_i) \in \mathbb{R} \times \mathbb{R}, i = 1, \dots, n$ from a model

$$y_i = f_0(x_i) + \varepsilon_i, i = 1, \dots, n$$

where $\varepsilon_i, i = 1, \dots, n$ are independent with $\mathbf{E}[\varepsilon_i] = 0$ and $\mathbf{E}[\varepsilon_i^2] = \sigma^2$. For simplicity, we treat the input points $x_i \in [0, 1], i = 1, \dots, n$ as fixed, satisfying the condition

$$\mathbf{P}_n(I) \geq c |I|, \text{ for any interval } I \subset [0, 1] \text{ with } |I| \geq 1/n,$$

where \mathbf{P}_n is the empirical distribution of the input points. Also assume that $x \in [0, 1]$. We also assume that the underlying regression function f_0 has a continuous, bounded derivative. Consider \widehat{f} , the kernel smoothing estimate with a boxcar kernel and bandwidth h .

1. Prove that the squared bias and variance of \widehat{f} , at an arbitrary point x , satisfy

$$\left(\mathbf{E} \left[\widehat{f}(x) \right] - f_0(x) \right)^2 \lesssim h^2 \text{ and } \mathbf{E} \left[\left(\widehat{f}(x) - \mathbf{E}[\widehat{f}(x)] \right)^2 \right] \lesssim \frac{1}{nh}$$

(Hint: use a Taylor expansion of f_0 around x .)

2. Derive the rate for the optimal choice of bandwidth h , and give the error rate for the corresponding kernel smoothing estimator.

b.

Consider the univariate k -th order local polynomial regression estimate, trained on the points $(x_i, y_i) \in \mathbb{R} \times \mathbb{R}, i = 1, \dots, n$, which we know can be expressed as

$$\widehat{f}(x) = \sum_{i=1}^n w_i(x) y_i,$$

for some weights $w_i(x), i = 1, \dots, n$.

1. Prove that at any point $x \in \mathbb{R}$, we have

$$\sum_{i=1}^n w_i(x) = 1 \quad \text{and} \quad \sum_{i=1}^n w_i(x)(x_i - x)^j = 0, \quad \text{for } j = 1, \dots, k.$$

2. Now assume the same model for the data as in part (a), and further assume that $x_i = 1/n, i = 1, \dots, n$ and $f_0 \in H_1(k+1, L)$ for a positive integer k and a constant $L > 0$. Also assume that $x \in [0, 1]$. Take \hat{f} to be the local polynomial regression estimate of order k , and compute the bias of \hat{f} at an arbitrary point x , using the results in the last part. (Hint: use a Taylor expansion, and the results of the last question. You can use the fact that $\sum_{i=1}^n |w_i| \leq C$ for some constant C that does not depend on n or h .) What do you conclude about the bias of local polynomial regression, compared that of kernel regression?
3. Why don't we just keep increasing the polynomial order k without end? You can answer this either with some theory, or a quick simulation. (Hint: consider the variance of \hat{f} . You can use the fact that $\sum_{i=1}^n w_i^2(x)$ is an increasing function of k .)

C.

Suppose that, in backfitting, we choose our univariate smoother just to be standard univariate linear regression. Prove that backfitting converges in one pass, and the resulting estimate is just standard multivariate linear regression.