

DATA 607 Statistical and Machine Learning

Session 2: Linear Smoothers

Matthew Greenberg

Department of Mathematics and Statistics
University of Calgary

26.02.2020

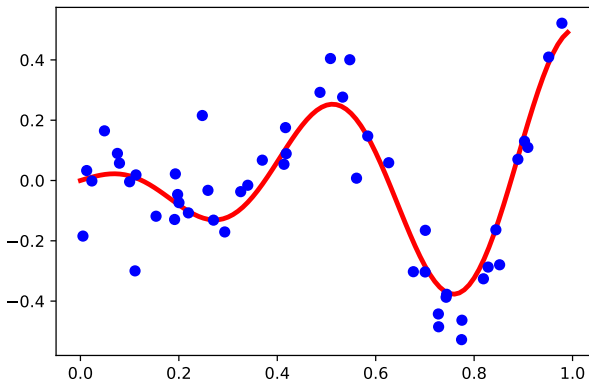
This Evening's Agenda

1 Linear Smoothers

- Definition
- Example: Linear Regression
- Example: k -Nearest Neighbors
- Example: Sliding Window

2 Evaluating Smoothers and Tuning Parameters

Smoothers are regression models for fitting “smooth” curves to data.



Linear Smoothers

Dataset:

$$(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)$$

Regression Model:

$$Y_i = r(\mathbf{x}_i) + \varepsilon_i$$

Definition

A *linear smoother* is an estimator of r of the form

$$\hat{r}(\mathbf{x}) = \sum_{i=1}^n w_i(\mathbf{x}) Y_i = \mathbf{w}(\mathbf{x}) \cdot \mathbf{Y}.$$

The $w_i(\mathbf{x})$ are called *weights*.

Linear smoother are so named because they are *linear functions of \mathbf{Y}* . Their graphs need not be lines!

Example 1: Linear Regression

Dataset:

$$(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n), \quad \mathbf{x}_i \in \mathbb{R}^{1 \times p}, \quad Y_i \in \mathbb{R}$$

Regression Model:

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \varepsilon \\ &= \underbrace{\begin{pmatrix} 1 & \mathbf{x}_i \end{pmatrix}}_{1 \times (p+1)} \boldsymbol{\beta} + \varepsilon, \end{aligned}$$

where

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^{(p+1) \times 1}.$$

Least-Squares Line:

$$\hat{r}(\mathbf{x}) = \mathbf{x}\hat{\beta}, \quad (*)$$

$\hat{\beta}$ is the least-squares solution of

$$X\beta = \mathbf{Y}, \quad \text{where} \quad X = \begin{pmatrix} 1 & \mathbf{x}_1 \\ \vdots & \vdots \\ 1 & \mathbf{x}_n \end{pmatrix} \in \mathbb{R}^{n \times (p+1)}.$$

Explicitly,

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y}.$$

Substituting into (*),

$$\hat{r}(\mathbf{x}) = \underbrace{(\mathbf{x}(X^T X)^{-1} X^T)}_{\mathbf{w}(\mathbf{x})} \mathbf{Y}.$$

Thus, \hat{r} is a linear smoother.

Coding Activity 1

- 1 Mock some data,

$$(x_0, y_0), \dots, (x_{19}, y_{19}), \quad x_i \in (0, 1), y_i \in \mathbb{R},$$

suitable for linear regression:

- 1 Choose β_0 and β_1 uniformly at random from the intervals $(0, 1)$ and $(-1, 0)$, respectively.
 - 2 Choose x_0, \dots, x_{19} uniformly at random from the interval $(-1, 1)$.
 - 3 Draw a random sample $\varepsilon_0, \dots, \varepsilon_{19}$ from $N(0, 0.2)$.
 - 4 Set $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$.
- 2 Plot y versus x . It should look roughly linear.
 - 3 Write a function that computes $\mathbf{w}(x)$ from x .
 - 4 Plot the graph of $y = \hat{r}(x)$.

Example 2: k -Nearest Neighbor Smoother

Data:

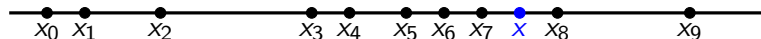
$$(\mathbf{x}_0, Y_0) \dots, (\mathbf{x}_{n-1}, Y_{n-1})$$

Regression Model:

$$Y_i = r(\mathbf{x}_i) + \varepsilon_i$$

k -Nearest Neighbors:

$N_k(\mathbf{x}) =$ the k elements of $\mathbf{x}_0, \dots, \mathbf{x}_{n-1}$ closest to \mathbf{x}



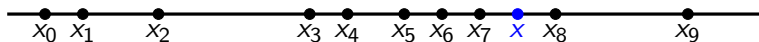
$$N_3(x) = \{x_6, x_7, x_8\}$$

Definition

The **k -nearest neighbor smoother** is the estimator of r is defined by

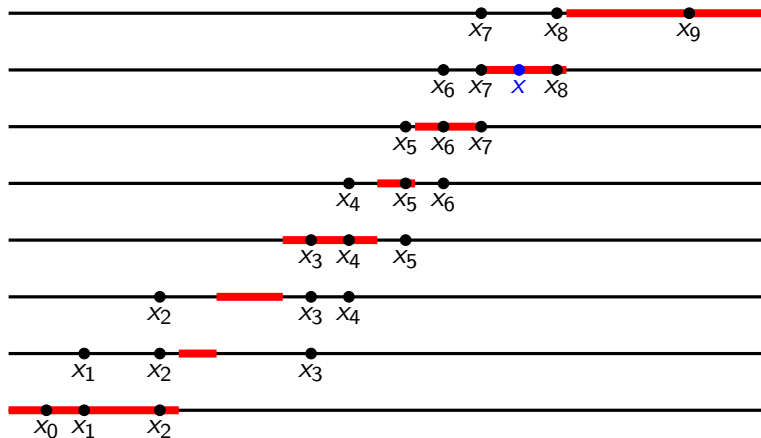
$$\hat{r}_k(\mathbf{x}) = \text{average of the } Y_i \text{ for which } \mathbf{x}_i \in N_k(\mathbf{x})$$

$$= \sum_{i=0}^{n-1} w_i(\mathbf{x}) Y_i, \quad \text{where} \quad w_i(\mathbf{x}) = \begin{cases} \frac{1}{k} & \text{if } \mathbf{x}_i \in N_k(\mathbf{x}), \\ 0 & \text{otherwise.} \end{cases}$$

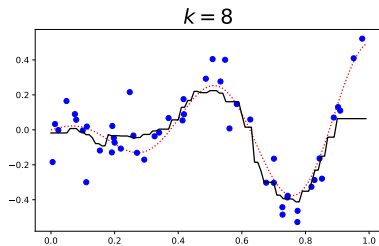
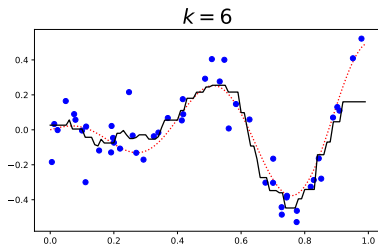
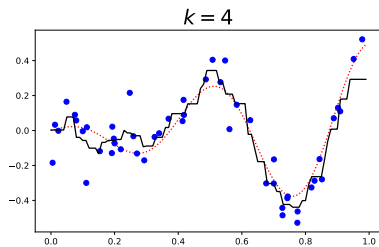
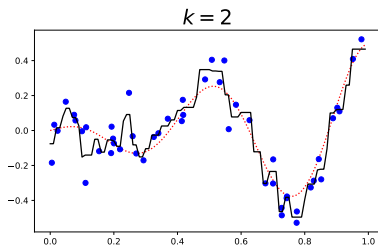


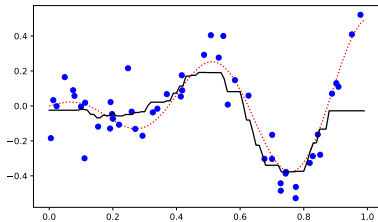
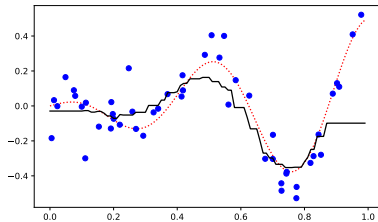
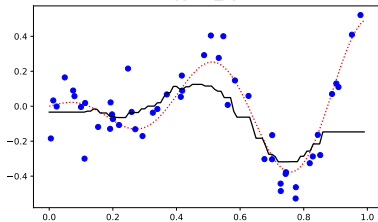
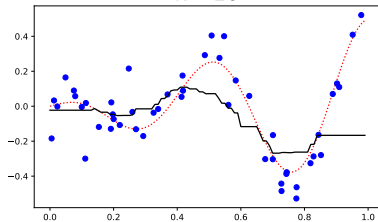
$$N_3(x) = \{x_6, x_7, x_8\}, \quad \hat{r}_3(x) = \frac{1}{3}(Y_6 + Y_7 + Y_8)$$

3-Nearest Neighbors:

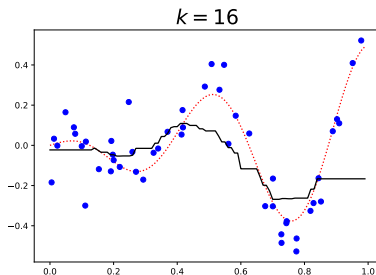
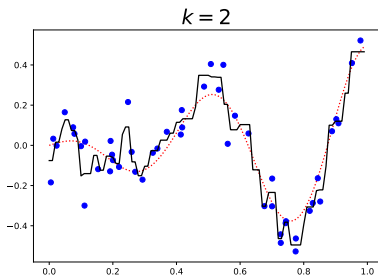


Sample Plots



$k = 10$  $k = 12$  $k = 14$  $k = 16$ 

Overfitting and Underfitting



For small k , the graphs of $\hat{r}(x)$ are too wiggly, fitting the noisy sample *too closely*. We've **overfit** (**undersmoothed**) the data.

For big k , the graphs of $\hat{r}(x)$ don't fit data very well (they're too flat). We've **underfit** (**oversmoothed**) the data.

Is there a “right” k ?

Example 3: The Sliding Window Smoother

Data:

$$(\mathbf{x}_0, Y_0) \dots, (\mathbf{x}_{n-1}, Y_{n-1})$$

Regression Model:

$$Y_i = r(\mathbf{x}_i) + \varepsilon_i$$

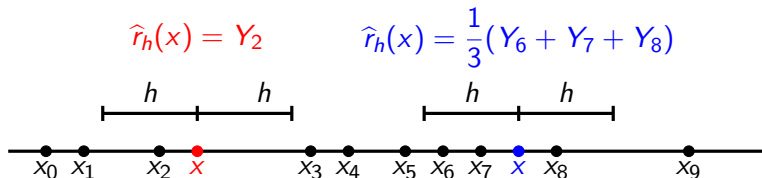
h -neighbors:

$$\begin{aligned} N_h(\mathbf{x}) &= \text{the elements of } \mathbf{x}_0, \dots, \mathbf{x}_{n-1} \text{ within a distance } h \text{ of } \mathbf{x} \\ &= \{\mathbf{x}_i : \|\mathbf{x}_i - \mathbf{x}\| < h\} \end{aligned}$$

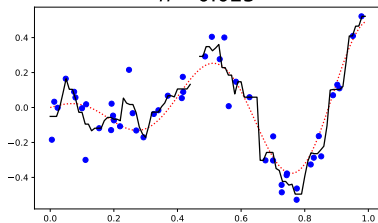
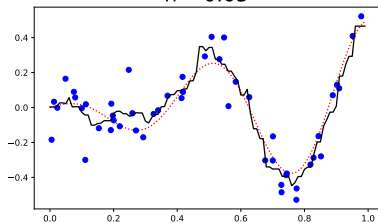
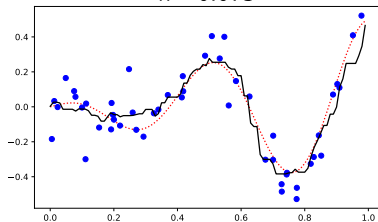
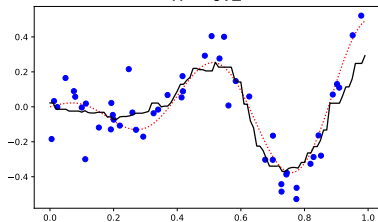
Definition

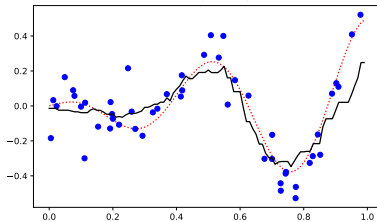
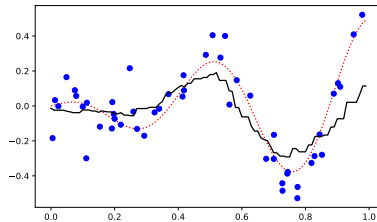
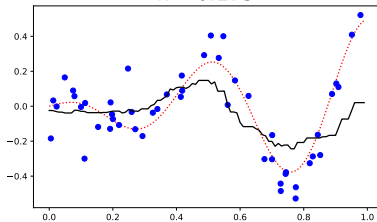
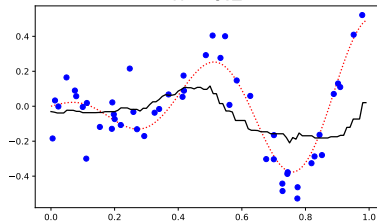
The **sliding window smoother** with bandwidth $h > 0$ is the estimator of r defined by

$$\begin{aligned}\hat{r}_h(\mathbf{x}) &= \text{average of the } Y_i \text{ for which } \|\mathbf{x} - \mathbf{x}_i\| < h \\ &= \frac{\sum_{\|\mathbf{x}_i - \mathbf{x}\| < h} Y_i}{\sum_{\|\mathbf{x}_i - \mathbf{x}\| < h} 1} \\ &= \sum_{i=0}^{n-1} w_i(\mathbf{x}) Y_i, \quad \text{where } w_i(\mathbf{x}) = \begin{cases} \frac{1}{\#N_h(\mathbf{x})} & \text{if } \mathbf{x}_i \in N_h(\mathbf{x}), \\ 0 & \text{otherwise.} \end{cases}\end{aligned}$$

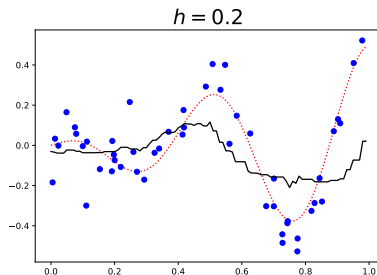
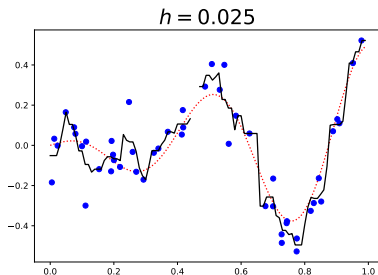


Sample Plots

 $h = 0.025$  $h = 0.05$  $h = 0.075$  $h = 0.1$ 

$h = 0.125$  $h = 0.15$  $h = 0.175$  $h = 0.2$ 

Overfitting and Underfitting



For small h , the graphs of $\hat{r}(x)$ are too wiggly, fitting the noisy sample *too closely*. We've **overfit** (**undersmoothed**) the data.

For large h , the graphs of $\hat{r}(x)$ don't fit data very well (they're too flat). We've **underfit** (**oversmoothed**) the data.

Is there a “right” h ?

Evaluating Smoothers: Predictive Accuracy

We evaluate a smoother based on the accuracy of its predictions.

Let $\hat{r}_{\mathcal{D}}$ be a smoother trained on the data set \mathcal{D} .

Assume we are in possession of a **test set**,

$$\mathcal{D}_{\text{test}} = \{(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)\},$$

i.e., a random sample drawn from the **same distribution** as — but **independently** of — \mathcal{D} .

The **expected prediction error** of $\hat{r}_{\mathcal{D}}$ is approximately the average error on $\mathcal{D}_{\text{test}}$:

$$\text{test error} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{r}_{\mathcal{D}}(\mathbf{x}_i))^2$$

Tuning Model Parameters Using a Test Set

Given a test set, $\mathcal{D}_{\text{test}}$, choose your parameters (k for k -nearest neighbors, h for sliding window) to **minimize**

$$\text{test error} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{r}_{\mathcal{D}}(\mathbf{x}_i))^2,$$

our proxy for the expected prediction error.

Compute the test error for a range of values of k or h , and choose the parameter value giving the smallest test error.

CODING ACTIVITY 2

But, we can work in an artificial situation where we *do* know $r(x)$:

- Choose:
 - a function, $r(x)$
 - x_1, \dots, x_n
 - a fine partition, P , of an interval containing the x_i
 - k (or h)
- for $k = 1, \dots, K$
 - for $j = 0, \dots, J - 1$:
 - Generate $y_i^{(j)}$, $1 \leq i \leq n$, by sampling from $N(r_k(x_i), \sigma^2)$.
 - Compute the losses $L(\hat{r}_k^{(j)}(x), r(x))$ at each x in P .
 - Average the losses over j to get the risks, $R(\hat{r}_k(x), r(x))$.
 - Average the risks over x to get the average risk, $R(\hat{r}_k, r)$.
- Plot $R(\hat{r}_k, r)$ vs k and choose k .

THE BIAS-VARIANCE DECOMPOSITION

Definition

The *bias* of $\hat{r}(x)$, as an estimator of $r(x)$, is

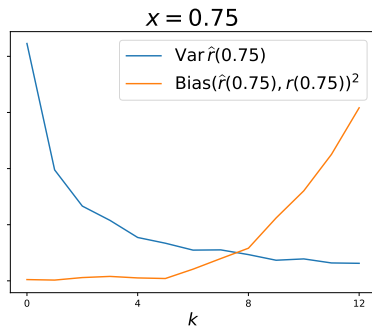
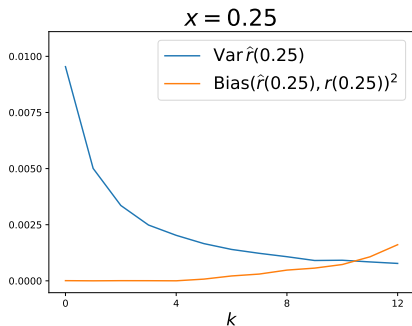
$$\text{Bias}(\hat{r}(x), r(x)) := \mathbb{E}[\hat{r}(x)] - r(x)$$

Bias-variance decomposition

$$R(\hat{r}(x), r(x)) = \text{Var} \hat{r}(x) + \text{Bias}(\hat{r}(x), r(x))^2$$

$\text{Var} \hat{r}(x)$ big / graph of \hat{r} wiggly / overfit / undersmoothed

$\text{Bias}(\hat{r}(x), r(x))$ big / graph of \hat{r} too flat / underfit / oversmoothed



Variance decreases as k increases.

Bias increases as k increases.

Define:

$$\text{Var } \hat{r} := \int R(\hat{r}(x), r(x)) dx, \quad \text{Bias}(\hat{r}, r)^2 := \int \text{Bias}(\hat{r}(x), r(x))^2 dx$$

Integrate the bias-variance decomposition:

$$R(\hat{r}, r) = \text{Var } \hat{r} + \text{Bias}(\hat{r}, r)^2$$

Since r is unknown, these quantities can only be estimated.

CODING ACTIVITY 3

- For our running synthetic example, plot

$$R(\hat{r}_k, r), \text{Var} \hat{r}_k, \text{ and } \text{Bias}(\hat{r}_k, r)^2$$

versus k on the same axes.

- Verify the integrated bias-variance decomposition using your computed $R(\hat{r}_k, r)$, $\text{Var} \hat{r}_k$, and $\text{Bias}(\hat{r}_k, r)^2$.

Reuse as much of your code from CODING ACTIVITY 2 as possible.

TRAINING ERROR

Definition

The *empirical risk* or *training error* is

$$\frac{1}{n} \sum_{i=1}^n L(\hat{r}(x_i), Y_i).$$

Empirical risk is not good estimator of R : We trained \hat{r} on the (x_i, Y_i) , making $L(\hat{r}(x_i), Y_i)$ biased downwards.

LEAVE ONE OUT CROSS-VALIDATION

Definition

The *leave one out cross validation (LOOCV)* score is

$$\hat{R}(h) := \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{r}_h^{(-i)}(x_i))^2,$$

where $\hat{r}_h^{(-i)}$ is computed using the subdataset of the original one by removing the point (x_i, Y_i) .

$\hat{R}(h)$ is typically a good estimator of the average risk of \hat{r}_h :

$$\hat{R}(h) \approx R(\hat{r}_h, r) + \sigma^2, \quad \text{where } \sigma^2 = \text{Var } Y_i.$$

We choose our smoothing parameter to be the one that minimizes $\hat{R}(h)$:

$$h := \underset{h}{\operatorname{argmin}} \hat{R}(h)$$

CODING ACTIVITY 4

For our running example, plot $\hat{R}(k)$ versus k . Compute the $\hat{R}(k)$ using `sklearn`'s `LeaveOneOut` class and the generator returned by its `split` method.

Which k that minimizes $\hat{R}(k)$? Compare with the results of CODING ACTIVITY 2.

For this optimal value of k , plot $\hat{r}_k^{(-i)}(x_i)$ versus x_i and $r(x)$ versus x on the same axes.

Remark: `LeaveOneOut` is fairly low-level; `sklearn` provides more convenient ways to tune hyperparameters using cross validation. Sometimes, though, it's necessary to work with the lower-level constructs.

K-FOLD CROSS VALIDATION

Partition $\{1, \dots, n\}$ into K folds, I_1, \dots, I_K , of roughly equal size. Let $\hat{r}_h^{(-I_j)}$ be the smoother computed using the subdataset of the original one obtained by removing (x_i, Y_i) , for $i \in I_j$.

Definition

The K -fold cross validation score is

$$\hat{R}_K(h) := \sum_{j=1}^K \frac{1}{|I_j|} \sum_{i \in I_j} (Y_i - r_h^{(-I_j)}(x_i))^2$$

We can use \hat{R}_K in place of \hat{R} to select a smoothing parameter. In practice, K is usually 5 (sklearn's default) or 10.

CODING ACTIVITY 5

Find the values of k that minimize the 3-, 5-, and 10-fold cross validation scores. Compute these scores using `cross_val_score` from `sklearn.model_selection`.

Confirm your results from CODING ACTIVITY 4 by performing LOOCV as n -fold cross validation, n being the size of the dataset.