

# 6. Kernel Smoothing Methods

오영민

# Index

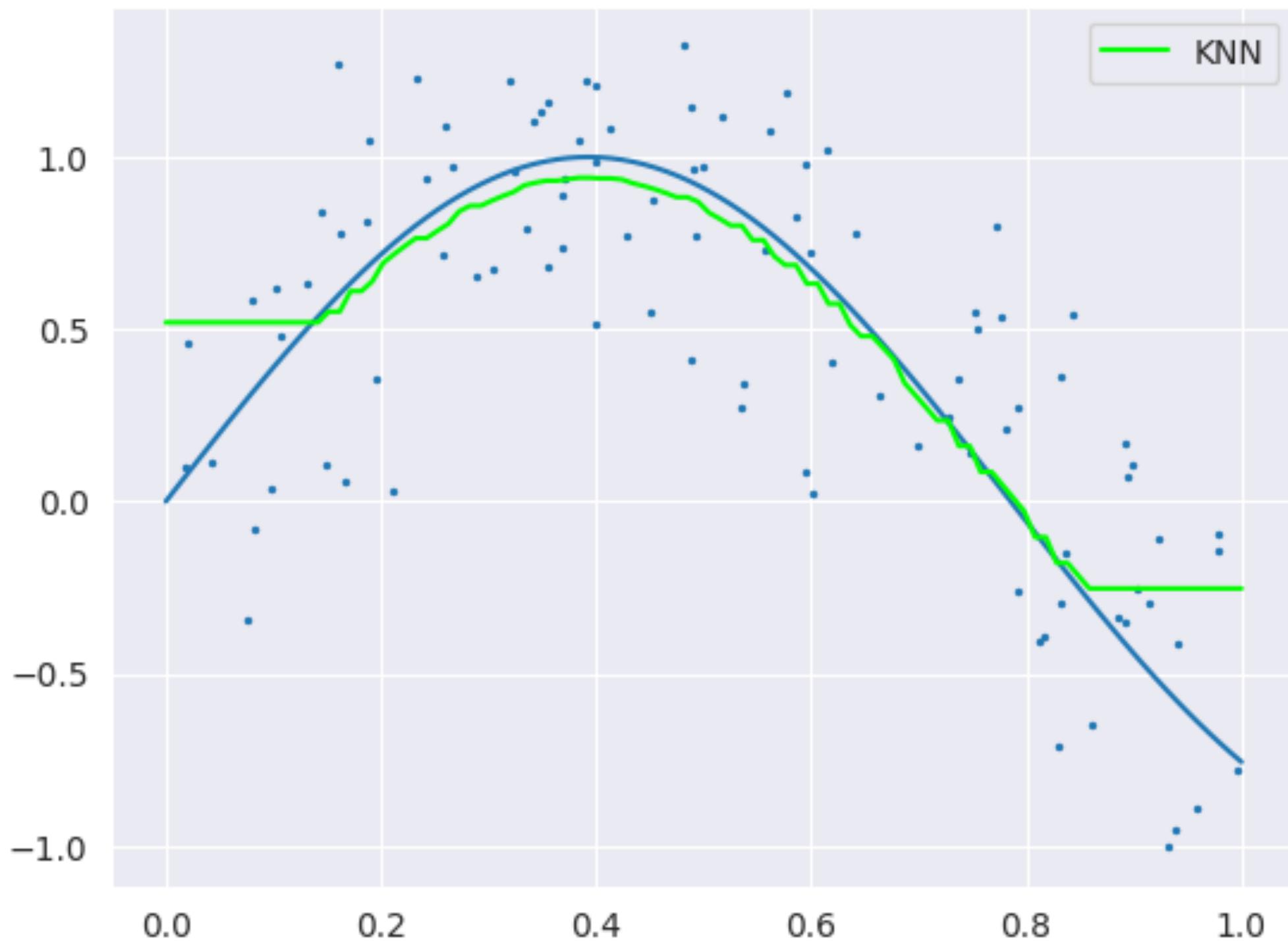
- Introduction
- One-Dimensional Kernel Smoothers
- Local Regression in  $\mathbb{R}^P$
- Local Likelihood and Other Models

# Introduction

- Using only those observations close to the target point  $x_0$  to fit the simple model s.t. resulting estimated  $\hat{f}$  is smooth in  $\mathbb{R}^p$  with weighting function of kernel  $K_\lambda(x_0 \cdot x_i)$ .
- This memory-based methods require in principle little or no training

# One-Dimensional Kernel Smoothers

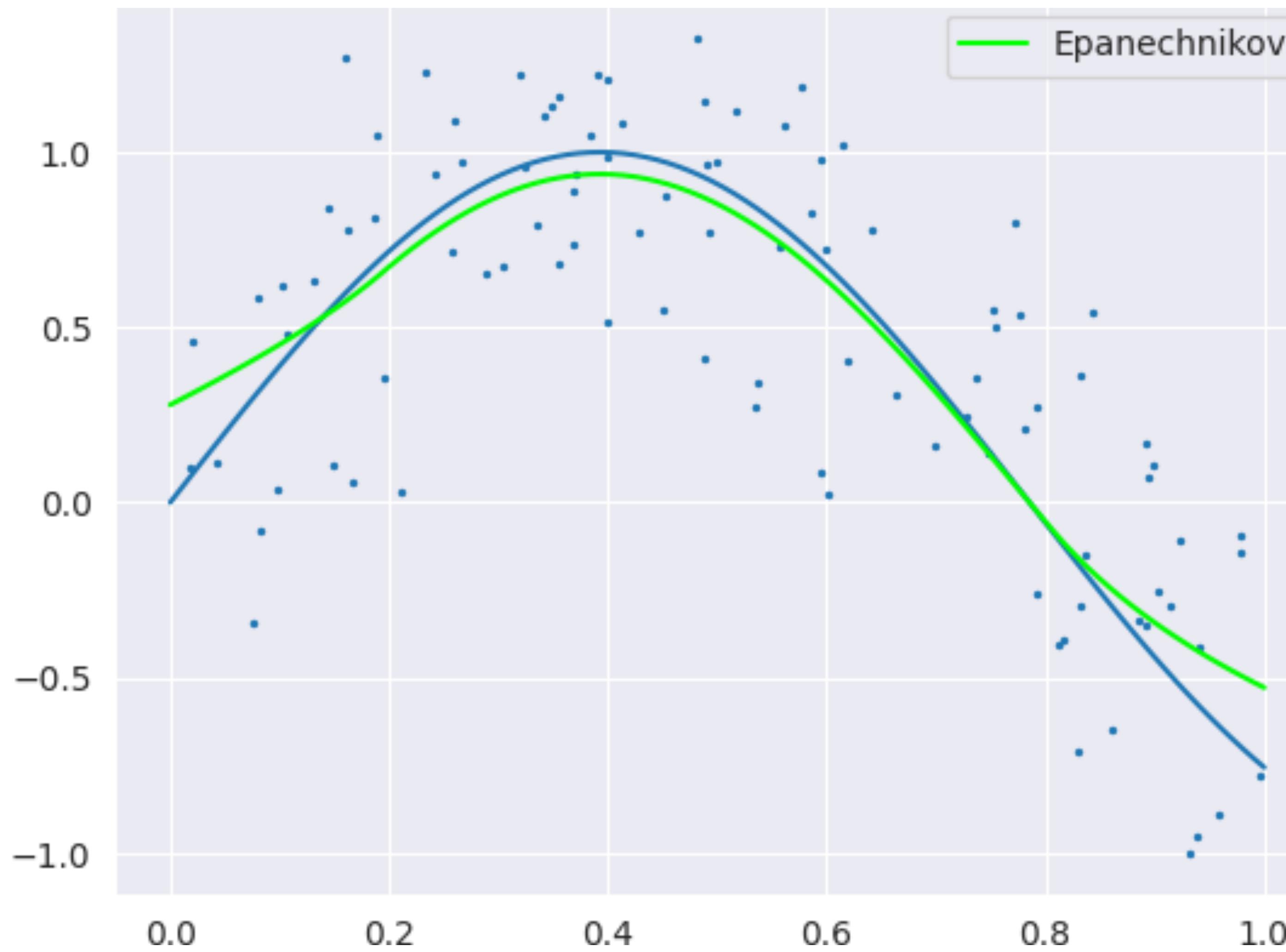
- **k-nearest-neighbor average:**  $\hat{f}(x) = \text{Ave}(y_i \mid x_i \in N_k(x))$  as an estimate of  $\mathbb{E}[Y \mid X = x]$
- Ex: In each panel 100 pairs  $Y_i = \sin(4X) + \epsilon$ ,  $X \sim U[0,1]$ ,  $\epsilon \sim \mathcal{N}(0,1/3)$



# One-Dimensional Kernel Smoothers

- **Nadaraya-Watson kernel-weighted average:**  $\hat{f}(x_0) = \frac{\sum_{i=1}^N K_\lambda(x_0, x_i) y_i}{\sum_{i=1}^N K_\lambda(x_0, x_i)}$
- **Kernel:**
  - Epanechnikov quadratic kernel:  $K_\lambda(x_0, x) = I[x \in (x_0 - \lambda, x_0 + \lambda)] \frac{3}{4} \left[ 1 - \left( \frac{x - x_0}{\lambda} \right)^2 \right]$  where  $\lambda$  is width or smoothing parameter. More general, use  $h_\lambda(x_0)$  instead  $\lambda$ .

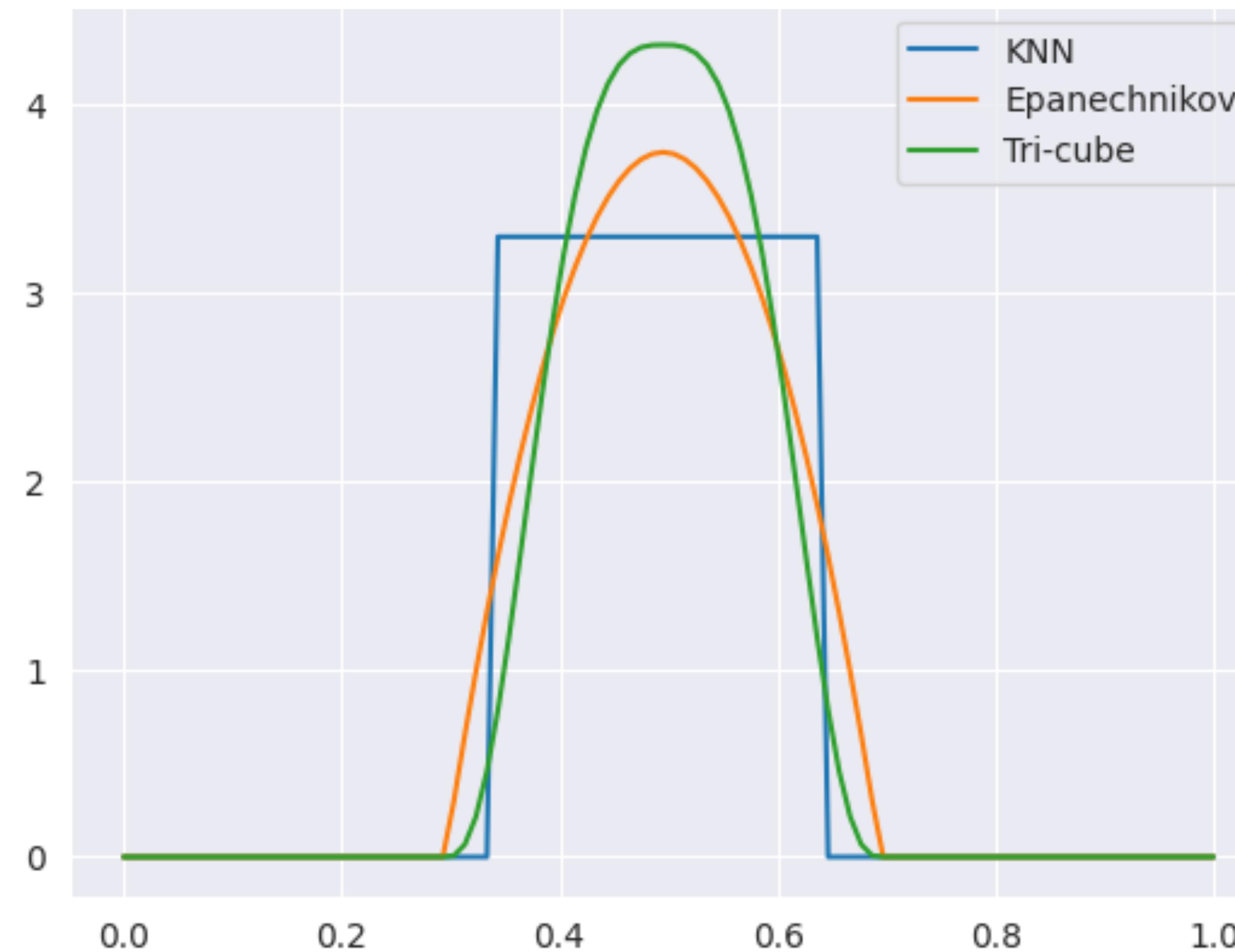
# One-Dimensional Kernel Smoothers



# One-Dimensional Kernel Smoothers

- **Tri-cube kernel:**  $K_\lambda(x_0, x) = I[x \in (x_0 - \lambda, x_0 + \lambda)] \left[ 1 - \left| \frac{x - x_0}{\lambda} \right|^3 \right]^3$
- **k-NN kernel:**  
$$K_k(x_0, x) = I[x \in N_k(x)] \frac{1}{2|x_0 - x_{[k]}|} = I[x \in (x_0 - |x_0 - x_{[k]}|, x_0 + |x_0 - x_{[k]}|)] \frac{1}{2|x_0 - x_{[k]}|}$$
where  $x_{[k]}$  is the k-th closest  $x_i$  to  $x_0$ .

# One-Dimensional Kernel Smoothers

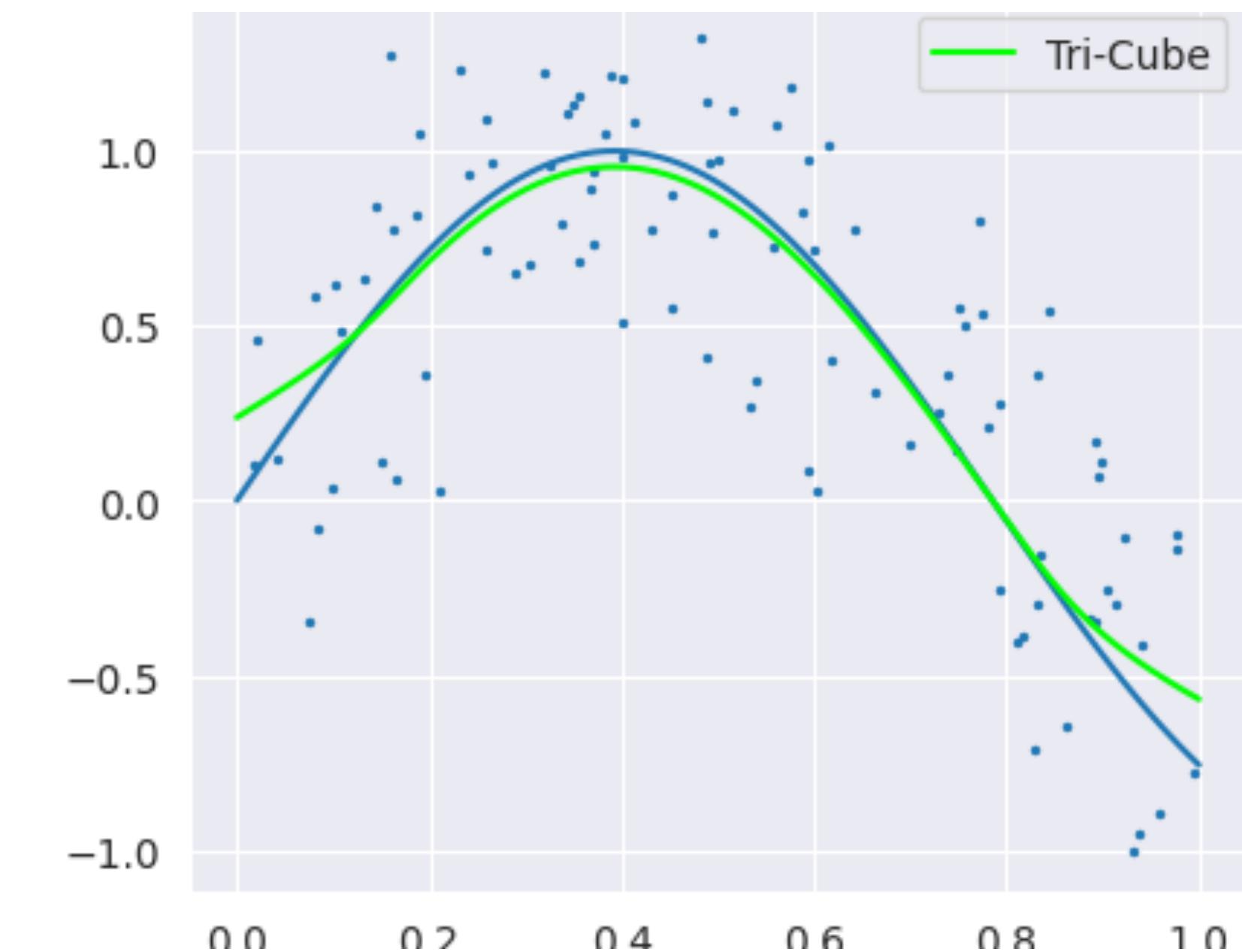
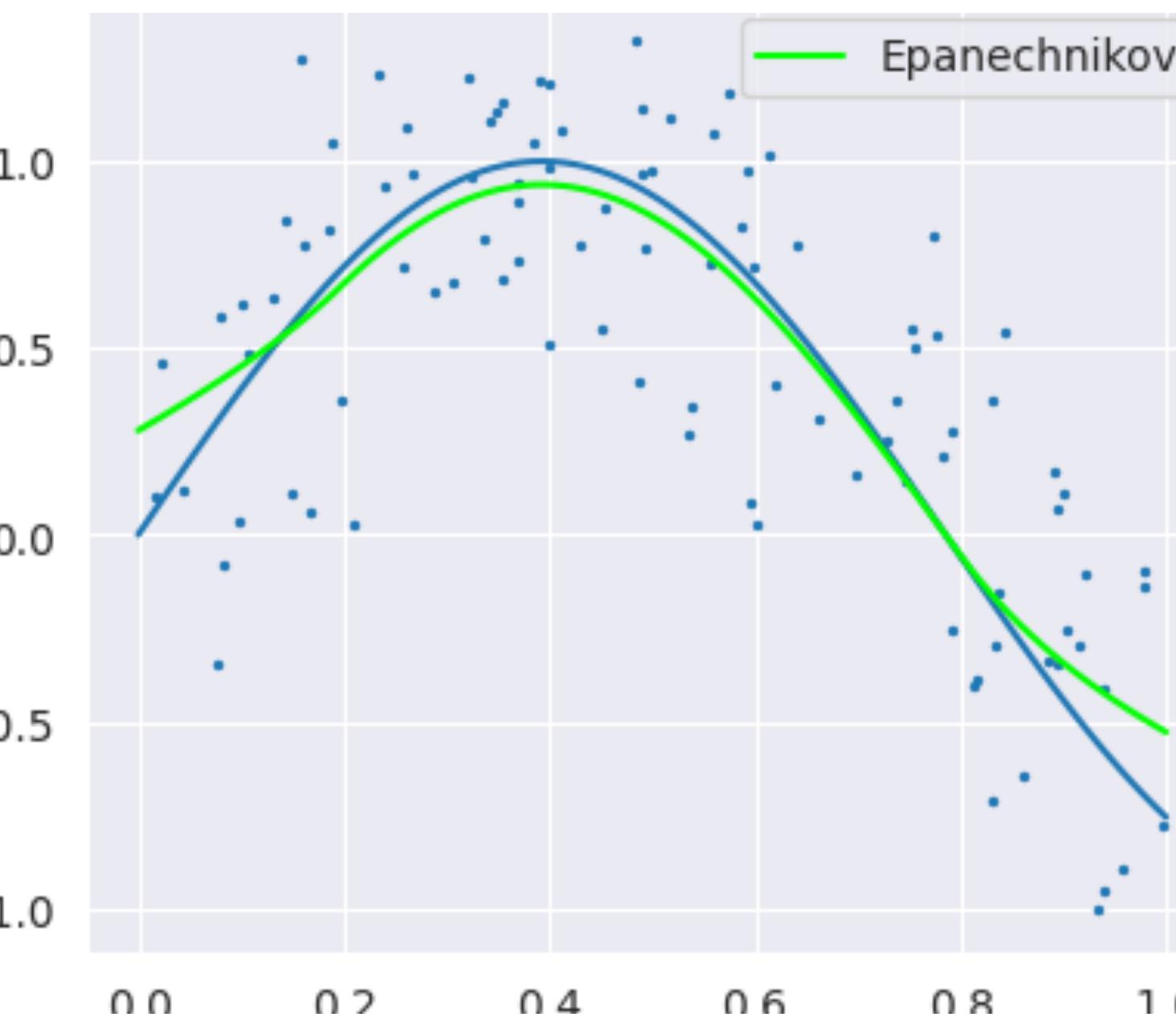
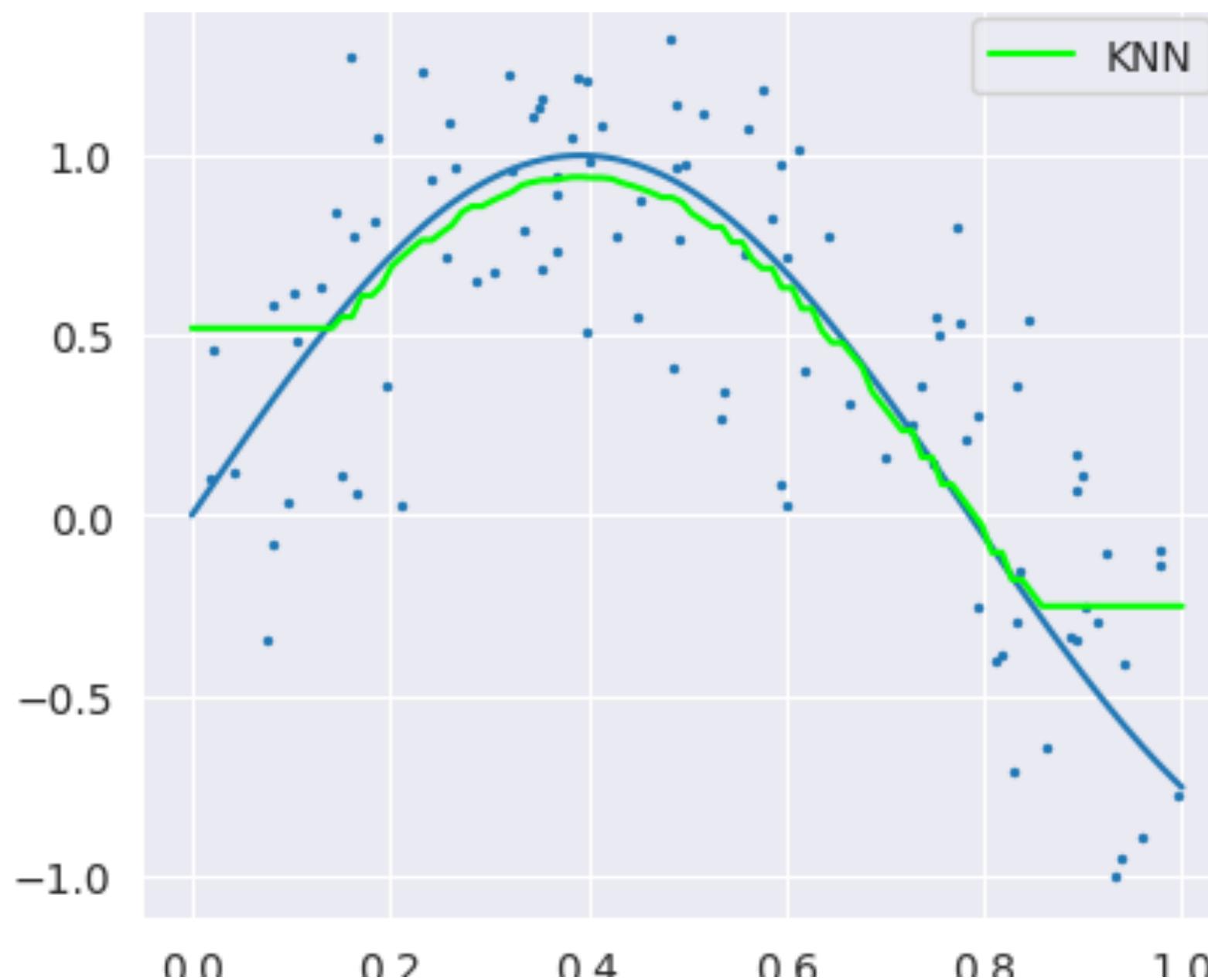


# One-Dimensional Kernel Smoothers

- The smoothing parameter  $\lambda$ , has to be determined. Large  $\lambda$  implies lower variance, but higher bias.
- NW kernel can be badly based on the **boundaries** of the domain, because of asymmetry of the kernel in that region.

# One-Dimensional Kernel Smoothers

- Most observations in neighborhood have higher mean than  $x_0$  s.t. near the boundaries
- By fitting straight lines rather then constant locally, we can remove this bias exactly to first order



# One-Dimensional Kernel Smoothers

- **Locally weighted regression:**  $\min_{\alpha(x_0)\beta(x_0)} \sum_{i=1}^N K_\lambda(x_0, x_i) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2$
- The estimate is then  $\hat{f}(x_0) = \hat{\alpha}(x_0) + \hat{\beta}(x_0)x_0$
- Notice that although we fit an entire linear model to the data in the region, we only use it to evaluate the fit at the single point  $x_0$ .

# One-Dimensional Kernel Smoothers

- Define the vector-valued function

$$b(x)^T = (1, x), \text{ and } B = \begin{bmatrix} b(x_1)^T \\ \vdots \\ b(x_N)^T \end{bmatrix}, W = \text{diag}(K_\lambda(x_0, x_1), \dots, K_\lambda(x_0, x_N))$$

$$\Rightarrow \text{Let } \beta = \begin{pmatrix} d_0 \\ \beta_0 \end{pmatrix} \in \mathbb{R}^2. \Rightarrow \sum_{n=1}^N k(x_0, x_n) [y_n - d_0 - \beta_0 - \beta_0^\top x_n]^2 = \|W(x_0)^{\frac{1}{2}}(y - B\beta)\|_2^2 \quad (\because k_\lambda(x_0, x_n) \geq 0 \text{ for all } n).$$

$$\Rightarrow L = \|W(x_0)^{\frac{1}{2}}(y - B\beta)\|_2^2 = (y - B\beta)^\top W(x_0)(y - B\beta) = y^\top W(x_0)y - y^\top W(x_0)B\beta - \beta^\top B^\top W(x_0)y + \beta^\top B^\top W(x_0)B\beta$$

$$\Rightarrow \frac{\partial L}{\partial \beta} = -y^\top W(x_0)B - y^\top W(x_0)B + \beta^\top B^\top W(x_0)B = 0 \Rightarrow \hat{\beta}^\top (B^\top W(x_0)B) = y^\top W(x_0)B \quad \hat{\beta}^\top = y^\top W(x_0)B(B^\top W(x_0)B)^{-1}$$

$$\therefore \hat{\beta} = (B^\top W(x_0)B)^{-1} B^\top W(x_0) y.$$

# One-Dimensional Kernel Smoothers

$$\begin{aligned}
 \hat{f}(x_0) &= b(x_0)^T (B^T W(x_0) B)^{-1} B^T W(x_0) y = b(x_0)^T \left( \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_N \end{pmatrix} \begin{pmatrix} k(x_0, x_1) & k(x_0, x_2) x_1 \\ \vdots & \vdots \\ k(x_0, x_N) & k(x_0, x_N) x_N \end{pmatrix} \right)^{-1} B^T W(x_0) y \\
 &= (1 \ x_0) \begin{pmatrix} \sum_n k(x_n, x_0) & \sum_n k(x_n, x_0) x_n \\ \sum_n k(x_n, x_0) x_n & \sum_n k(x_n, x_0) x_n^2 \end{pmatrix}^{-1} \begin{pmatrix} k(x_0, x_1) & \dots & k(x_0, x_N) \\ k(x_1, x_0) x_1 & \dots & k(x_N, x_0) x_N \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \\
 &= (1 \ x_0) \frac{1}{\sum_n k(x_n, x_0) \sum_n k(x_n, x_0) x_n^2 - (\sum_n k(x_n, x_0) x_n)^2} \begin{pmatrix} \sum_n k(x_n, x_0) x_n^2 & -\sum_n k(x_n, x_0) x_n \\ -\sum_n k(x_n, x_0) x_n & \sum_n k(x_n, x_0) \end{pmatrix} \begin{pmatrix} \sum_n k(x_n, x_0) y_n \\ \sum_n k(x_n, x_0) x_n y_n \end{pmatrix} \\
 &= \frac{1}{\sum_n k(x_n, x_0) \sum_n k(x_n, x_0) x_n^2 - (\sum_n k(x_n, x_0) x_n)^2} (\sum_n k(x_n, x_0) x_n^2 - x_0 \sum_n k(x_n, x_0) x_n - \sum_n k(x_n, x_0) x_n + x_0 \sum_n k(x_n, x_0)) \begin{pmatrix} \sum_n k(x_n, x_0) y_n \\ \sum_n k(x_n, x_0) x_n y_n \end{pmatrix} \\
 &= \frac{1}{\sum_n k(x_n, x_0) \sum_n k(x_n, x_0) x_n^2 - (\sum_n k(x_n, x_0) x_n)^2} (\sum_n k(x_n, x_0) x_n^2 - x_0 \sum_n k(x_n, x_0) x_n) \sum_n k(x_n, x_0) y_n + (x_0 \sum_n k(x_n, x_0) - \sum_n k(x_n, x_0) x_n) \sum_n k(x_n, x_0) x_n y_n \\
 &= \sum_{n=1}^N b_n(x_0) y_n. \text{ where } b_n(x_0) \text{ does not involve } y_n \Rightarrow \hat{f}(x_0) \text{ is linear in } y_0.
 \end{aligned}$$

# One-Dimensional Kernel Smoothers

- Thus,  $\hat{f}(x_0) = \sum_{i=1}^N l_i(x_0) y_i$  where 
$$l_i(x_0) = \frac{\left[ \sum_j k(d_j, x_0) x_j - d_0 \sum_j k(d_j, x_0) d_j \right] K(d_n, x_0) + \left[ d_0 \sum_j k(d_j, x_0) - \sum_j k(d_j, x_0) d_j \right] K(d_n, x_0) d_n}{\sum_j k(d_j, x_0) \sum_j k(d_j, x_0) x_j' - (\sum_j k(d_j, x_0) x_j)'}$$
- The estimate is linear in the  $y_i$ . These weights  $l_i(x_0)$  are sometimes referred to as the equivalent kernel.
- $\sum_{i=1}^N l_i(x_0) = 1$  and  $\sum_{i=1}^N (x_i - x_0) l_i(x_0) = 0$

# One-Dimensional Kernel Smoothers

$$\Rightarrow \text{def } \sum_n k(d_n, d_0) = \alpha, \quad \sum_n k(d_n, d_0) d_n = \beta, \quad \sum_n k(d_n, d_0) d_n^2 = \tau$$

$$\Rightarrow \sum_{n=1}^N d_n k(d_0) = \frac{\alpha r - d_0 \cancel{\alpha} + \cancel{\lambda_0} \cancel{\beta} - \beta^2}{\alpha r - \beta^2} = \frac{\alpha r - \beta^2}{\alpha r - \beta^2} = 1.$$

$$\sum_{n=1}^N (d_n - d_0) k_n(d_0) = \frac{1}{\alpha r - \beta^2} \left[ (r - \cancel{\lambda_0} \beta) \underbrace{\sum_n k(d_n, d_0) (d_n - d_0)}_{= \beta - \cancel{\lambda_0} \alpha} + (\cancel{\lambda_0} \alpha - \beta) \underbrace{\sum_n k(d_n, d_0) d_n (d_n - d_0)}_{= \tau - \cancel{\lambda_0} \beta} \right]$$

$$= \frac{1}{\alpha r - \beta^2} \left[ \cancel{\beta r} - \cancel{\lambda_0} \cancel{\alpha r} - \cancel{\lambda_0} \beta^2 + \cancel{\lambda_0} \cancel{\alpha} \beta + \cancel{\lambda_0} \cancel{\alpha r} - \cancel{\lambda_0} \beta^2 - \cancel{\beta r} + \cancel{\lambda_0} \beta^2 \right] = 0.$$

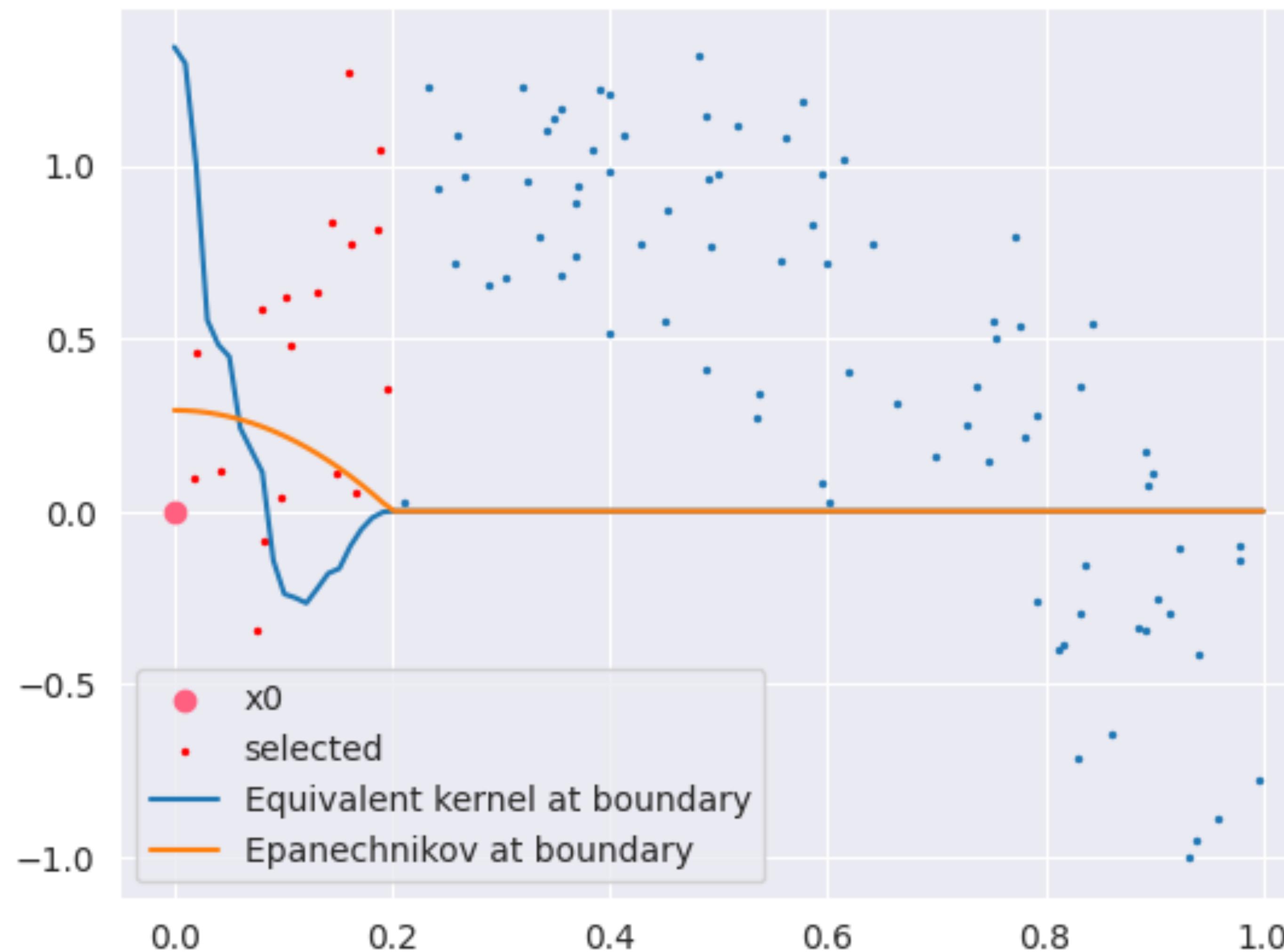
# One-Dimensional Kernel Smoothers

- **(Automatic kernel carpentry)** Local linear regression automatically modifies the kernel to correct the bias exactly to first order.
- Let  $f$  be true function,  $Y_i = f(x_i) + \epsilon_i$  where  $\epsilon_i$  are i.i.d with  $\mathbb{E}[\epsilon_i] = 0, \text{Var}[\epsilon_i] = \sigma^2$ .
- Then by Taylor series expansion,
$$\begin{aligned}\mathbb{E}\hat{f}(x_0) &= \sum_{i=1}^N l_i(x_0)f(x_i) \\ &= f(x_0) \sum_{i=1}^N l_i(x_0) + f'(x_0) \sum_{i=1}^N (x_i - x_0)l_i(x_0) \\ &\quad + \frac{f''(x_0)}{2} \sum_{i=1}^N (x_i - x_0)^2 l_i(x_0) + R,\end{aligned}$$

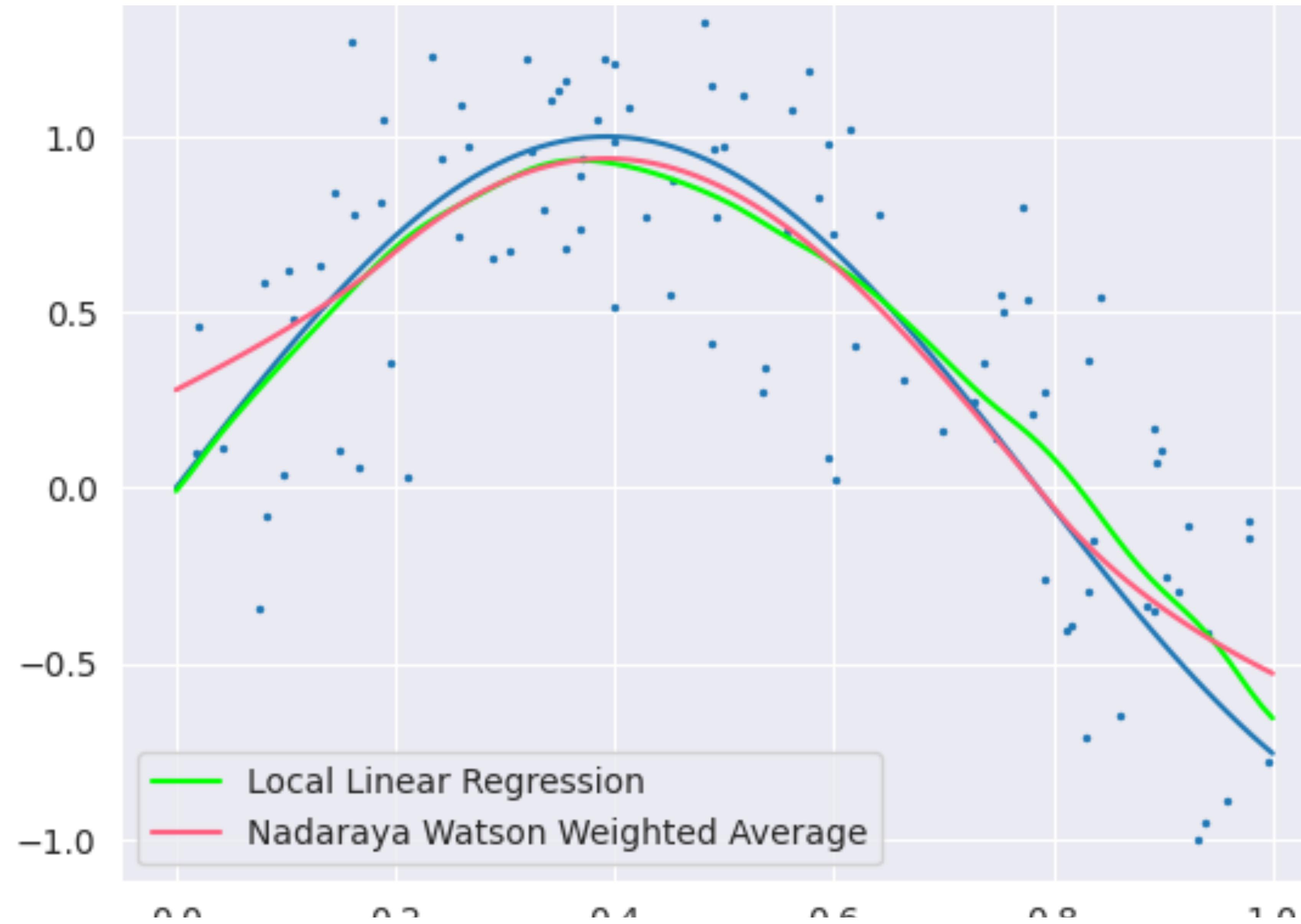
# One-Dimensional Kernel Smoothers

- Thus, the bias is  $\mathbb{E}\hat{f}(x_0) - f(x_0) = \frac{f''(x_0)}{2} \sum_i (x_i - x_0)^2 l_i(x_0) + R$  which depends only on quadratic and higher-order terms in the expansion of  $f$ .
- Under suitable smoothness assumptions, remained term is typically small.
- **(Trimming the hills)** Local linear fits tend to be biased in regions of curvature of the true function. -> local quadratic regression is generally able to correct this bias.

# One-Dimensional Kernel Smoothers



# One-Dimensional Kernel Smoothers



# One-Dimensional Kernel Smoothers

- Local Polynomial Regression:

$$\min_{\alpha(x_0), \beta_j(x_0), j=1, \dots, d} \sum_{i=1}^N K_\lambda(x_0, x_i) \left[ y_i - \alpha(x_0) - \sum_{j=1}^d \beta_j(x_0) x_i^j \right]^2 \text{ with solution}$$

$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \sum_{j=1}^d \hat{\beta}_j(x_0) x_0^j$$

# One-Dimensional Kernel Smoothers

Thm. Let  $b_j(t_0) = \sum_{n=1}^N (x_n - t_0)^j l_n(t_0)$  for local polynomial regression of degree  $k$ .

Then  $b_j(t_0) = \begin{cases} 1 & j=0 \\ 0 & j=1, \dots, k \end{cases}$

pf)  $B = \begin{pmatrix} 1 & x_1 & x_1^k \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^k \end{pmatrix} = \begin{pmatrix} b(x_1)^T \\ \vdots \\ b(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times (k+1)}$ ,  $w(t_0) = \text{diag}(k_1(t_0, x_1) \dots k_1(t_0, x_N))$

$\Rightarrow$  let  $\beta = (\alpha(t_0) \ \beta_1(t_0) \ \dots \ \beta_k(t_0))^T \in \mathbb{R}^{k+1}$

$$\Rightarrow \sum_{n=1}^N k_n(t_0, x_n) [y_n - \alpha(t_0) - \beta_1(t_0)x_1 - \dots - \beta_k(t_0)x_n^k]^2 = \|w(t_0)^{\frac{1}{2}}(y - B\beta)\|_2^2 \Rightarrow \hat{\beta} = (B^T w(t_0) B)^{-1} B^T w(t_0) y$$

$$\Rightarrow l(t_0)^T = b(t_0)^T (B^T w(t_0) B)^{-1} B^T w(t_0) \in \mathbb{R}^{1 \times N}. \text{ Note that } l(t_0)^T B = b(t_0)^T$$

$$\Rightarrow l(t_0)^T B = (l_1(t_0) \dots l_N(t_0)) \begin{pmatrix} 1 & x_1 & \dots & x_1^k \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \dots & x_N^k \end{pmatrix} = (\sum_n l_n(t_0) \quad \sum_n l_n(t_0)x_1 \quad \dots \quad \sum_n l_n(t_0)x_N^k) = (1 \ x_0 \ \dots \ x_0^k) \Rightarrow b_0(t_0) = \sum_n l_n(t_0) = 1.$$

$$\Rightarrow \sum_{n=1}^N (x_n - t_0)^j l_n(t_0) = \sum_{n=1}^N \sum_{k=0}^j \binom{j}{k} x_n^k (-t_0)^{j-k} l_n(t_0) = \sum_{k=0}^j \binom{j}{k} (-t_0)^{j-k} \sum_{n=1}^N x_n^k l_n(t_0) = \sum_{k=0}^j \binom{j}{k} (-t_0)^{j-k} x_0^k = (t_0 - x_0)^j = 0. \text{ for } j=1, \dots, k.$$

# One-Dimensional Kernel Smoothers

Corollary. (Relation between degrees & bias).  $E[\hat{f}(x_0)] - f(x_0) = R$  where  $R$  involves  $(k+1)$ - and higher-order derivatives of  $f$ .

$$\text{pf)} E[\hat{f}(x_0)] - f(x_0) = \sum_{n=1}^N l_n(x_0) f(x_n) - f(x_0) = \left[ f(x_0) \underbrace{\sum_{n=1}^N l_n(x_0)}_{=1} + f'(x_0) \sum_{n=1}^N l_n(x_0)(x_n - x_0) + \dots + \frac{f^{(k)}(x_0)}{k!} \sum_{n=1}^N l_n(x_0)(x_n - x_0)^k + R \right] - f(x_0)$$

$$= R.$$

- **(Bias-Variance Tradeoff)** There is of course a price to be paid for bias reduction, and that is increased variance

$$\text{• } \text{Var} [\hat{f}(x_0)] = \sum_{i=1}^N l_i(x_0)^2 \sigma^2 = \|l(x_0)\|_2^2 \text{ where } l(x_0) = (l_1(x_0), \dots, l_N(x_0))^T$$

# One-Dimensional Kernel Smoothers

Thm. (Relation between degrees & variance).  $\|\ell(x_0)\| = \|\ell_1(x_0), \dots, \ell_N(x_0)\|$  increases with  $d$ .

$$\begin{aligned} \|\ell(x_0)\|_2^2 &= [b(x_0)^T (B^T W(x_0) B)^{-1} B^T W(x_0)] [W(x_0) B (B^T W(x_0) B)^{-1} b(x_0)] \\ &= b(x_0)^T (B^T W(x_0) B)^{-1} B^T W(x_0) B (B^T W(x_0) B)^{-1} b(x_0) \\ &= b(x_0)^T B^{-1} (B^T)^{-1} b(x_0) = b(x_0)^T (B^T B)^{-1} b(x_0) \end{aligned}$$

$$\Rightarrow \text{def } b'(x_0) = (b(x_0) \ x_0^{d+1}), \quad B = (B \ b) \text{ where } b = (x_1^{d+1} \ \dots \ x_N^{d+1})^T \in \mathbb{R}^{d+1}$$

$$\begin{aligned} \Rightarrow \|\ell'(x_0)\|_2^2 &= \begin{pmatrix} b(x_0)^T \\ x_0^{d+1} \end{pmatrix} \left( \begin{pmatrix} B^T \\ b^T \end{pmatrix} (B \ b) \right)^{-1} \begin{pmatrix} b(x_0) \ x_0^{d+1} \end{pmatrix} \Rightarrow \text{check that } \|\ell'(x_0)\|_2 \geq \|\ell(x_0)\|_2 \\ &= \begin{pmatrix} b(x_0)^T \\ x_0^{d+1} \end{pmatrix} \begin{pmatrix} B^T B & B^T b \\ B^T b & \underbrace{b^T b}_{\in \mathbb{R}} \end{pmatrix}^{-1} \begin{pmatrix} b(x_0) \ x_0^{d+1} \end{pmatrix} \end{aligned}$$

# One-Dimensional Kernel Smoothers

$$\begin{pmatrix} B^T B & B^T b \\ b^T B & b^T b \end{pmatrix} = \begin{pmatrix} A & B \\ B^T & C \end{pmatrix} \text{ form.}$$

$$\begin{pmatrix} I & 0 \\ -B^T A^{-1} & I \end{pmatrix} \begin{pmatrix} A & B \\ B^T & C \end{pmatrix} \begin{pmatrix} I - A^{-1}B \\ 0 & I \end{pmatrix} = \begin{pmatrix} A & 0 \\ 0 & C - B^T A^{-1}B \end{pmatrix}$$

$$\left\langle \begin{pmatrix} A & B \\ B^T & C \end{pmatrix}^{-1} = \begin{pmatrix} I & -A^T B \\ 0 & I \end{pmatrix} \begin{pmatrix} A^{-1} & 0 \\ 0 & (C - B^T A^{-1}B)^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ -B^T A^{-1} & I \end{pmatrix} = \begin{pmatrix} A^{-1} + A^{-1}B(C - B^T A^{-1}B)^{-1}B^T A^{-1} \\ -(C - B^T A^{-1}B)^{-1}B^T A^{-1} \end{pmatrix} \right. \\ \left. - A^T B(C - B^T A^{-1}B)^{-1} \right\rangle \\ (C - B^T A^{-1}B)^{-1}$$

$$\Rightarrow \begin{pmatrix} B^T B & B^T b \\ b^T B & b^T b \end{pmatrix} = \begin{pmatrix} (B^T B)^{-1} + (B^T B)^{-1}B^T b(B^T b - b^T B(B^T B)^{-1}B^T b)^{-1}b^T B(B^T B)^{-1} & - (B^T B)^{-1}B^T b(B^T b - b^T B(B^T B)^{-1}B^T b)^{-1} \\ -(B^T b - b^T B(B^T B)^{-1}B^T b)^{-1}b^T B(B^T B)^{-1} & (B^T b - b^T B(B^T B)^{-1}B^T b)^{-1} \end{pmatrix} \\ = \begin{pmatrix} (B^T B)^{-1} + \alpha \beta \beta^T & -\alpha \beta \\ -\alpha \beta^T & \alpha \end{pmatrix}$$

$\beta \in \mathbb{R}^N$   
 $\alpha \in \mathbb{R}$

# One-Dimensional Kernel Smoothers

$$\begin{aligned}\Rightarrow \|\ell'(\lambda_0)\|_2^2 &= b(\lambda_0)^T [ (B^T B)^{-1} + \alpha \beta \beta^T ] b(\lambda_0) - b(\lambda_0)^T \alpha \beta \lambda_0^{d+1} - \lambda_0^{d+1} \alpha \beta^T b(\lambda_0)^T + \alpha (\lambda_0^{d+1})^2 \\ &= \|\ell(\lambda_0)\|_2^2 + \alpha \left\| \lambda_0^{d+1} - b(\lambda_0)^T \beta \right\|_2^2 \\ &\quad \text{Projection matrix.} \\ \Rightarrow \text{Enough to show that } b^T b - b^T B (B^T B)^{-1} B^T b &\geq 0 \Rightarrow b^T b - b^T B (B^T B)^{-1} B^T b = b^T (I - B (B^T B)^{-1} B^T) b \geq 0 \quad \square.\end{aligned}$$

- Selecting the parameter  $d, \lambda \rightarrow$  Cross validation
- Ex. Leave-One-Out Cross-Validation

# One-Dimensional Kernel Smoothers

In Local regression, using Smoother matrix,

$$\hat{f} = B(B^T W(x_0) B)^{-1} B^T W(x_0) y = \begin{bmatrix} b(x_1)^T \\ \vdots \\ b(x_N)^T \end{bmatrix} (B^T W(x_0) B)^{-1} B^T W(x_0) y = \begin{bmatrix} l(x_1)^T \\ \vdots \\ l(x_N)^T \end{bmatrix} y = \begin{bmatrix} l_1(x_1) & \dots & l_N(x_1) \\ \vdots & \ddots & \vdots \\ l_1(x_N) & \dots & l_N(x_N) \end{bmatrix} y = S_\lambda y.$$

• LOOCV:  $\text{CV}(\hat{f}_\lambda) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}^{-i}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N \left( \frac{y_i - \hat{f}(x_i)}{1 - S_{\lambda ii}} \right)^2$

# One-Dimensional Kernel Smoothers

$$\text{If } \{S_\lambda\}_{n,n} = b(t_n)^T (B^T W(t_0) B)^{-1} b(t_n) k_\lambda(t_0, t_n)$$

$$\Rightarrow \text{Let } B_{-i}, Y_{-i}, W(t_0)_{-i} \text{ be removed } b(t_n), y_n, k_\lambda(t_0, t_n) \text{ i.e. } B_{-i} = \begin{bmatrix} b(t_1)^T \\ \vdots \\ 0^T \\ b(t_N)^T \end{bmatrix}, Y_{-i} = \begin{bmatrix} y_1 \\ \vdots \\ 0 \\ y_N \end{bmatrix}.$$

$$\Rightarrow \hat{f}^{-i}(t_0) = b(t_0)^T (B_{-i}^T W(t_0)_{-i} B_{-i})^{-1} B_{-i}^T W(t_0)_{-i} Y_{-i}$$

$$\Rightarrow B_{-i}^T W(t_0)_{-i} B_{-i} = \sum_{j=1}^N b(t_j) k_\lambda(t_0, t_j) b(t_j)^T - b(t_0) k_\lambda(t_0, t_n) b(t_n)^T$$

$$\left\langle B_{-i}^T W(t_0)_{-i} Y_{-i} = \sum_{j=1}^N b(t_j) k_\lambda(t_0, t_j) y_j - b(t_0) k_\lambda(t_0, t_n) y_n \right\rangle$$

$\Rightarrow$  By Woodbury matrix identity,

$$(B_{-i}^T W(t_0)_{-i} B_{-i})^{-1} = (B^T W(t_0) B - b(t_n) k_\lambda(t_0, t_n) b(t_n)^T)^{-1}$$

$$= (B^T W(t_0) B)^{-1} + \frac{1}{\frac{1}{k_\lambda(t_0, t_n)} - b(t_n)^T (B^T W(t_0) B)^{-1} b(t_n)} \cdot (B^T W(t_0) B)^{-1} b(t_n) b(t_n)^T (B^T W(t_0) B)^{-1}.$$

$$= (B^T W(t_0) B)^{-1} + \frac{k_\lambda(t_0, t_n)}{1 - \{S_\lambda\}_{n,n}} \cdot (B^T W(t_0) B)^{-1} b(t_n) b(t_n)^T (B^T W(t_0) B)^{-1}.$$

# One-Dimensional Kernel Smoothers

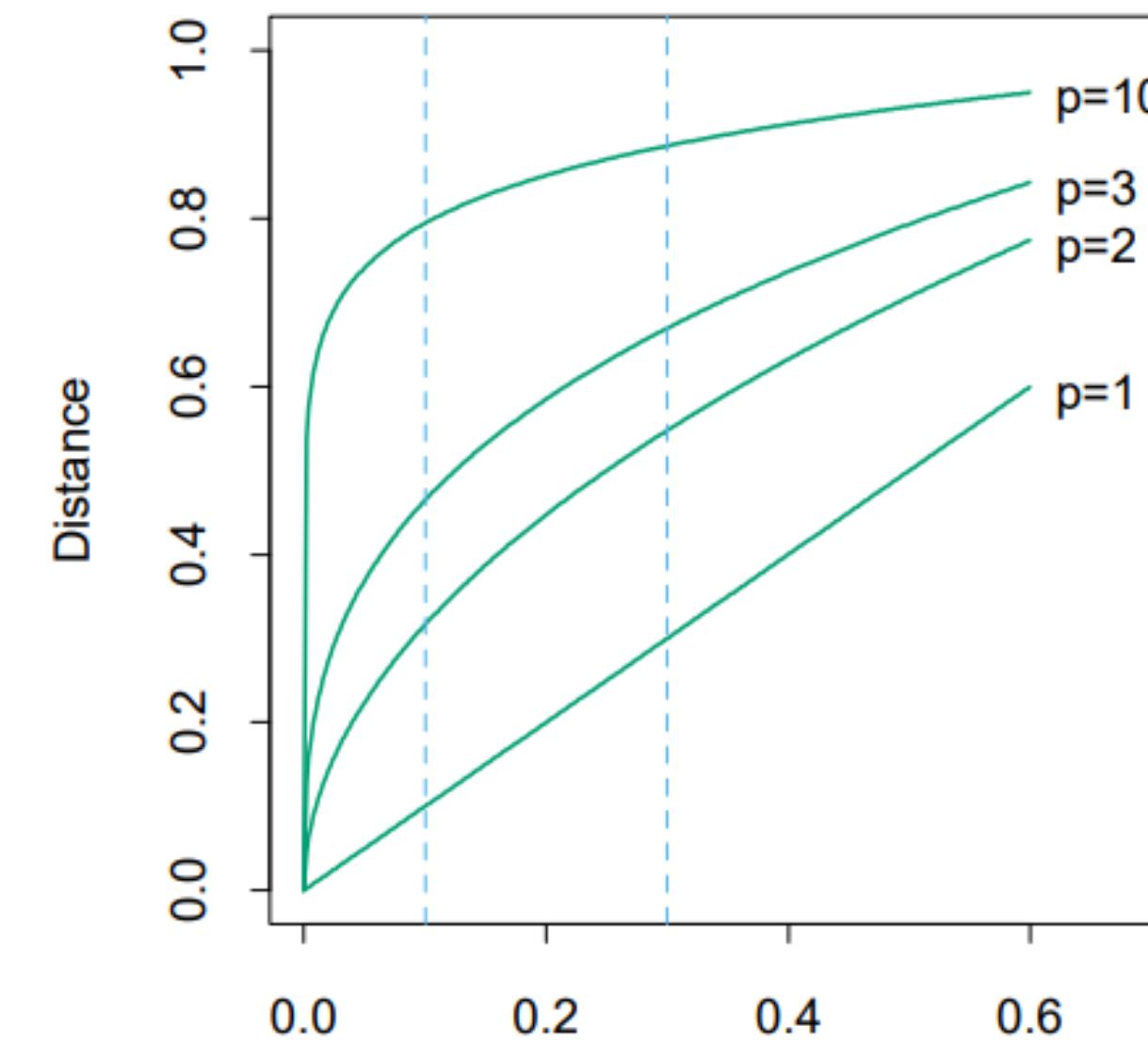
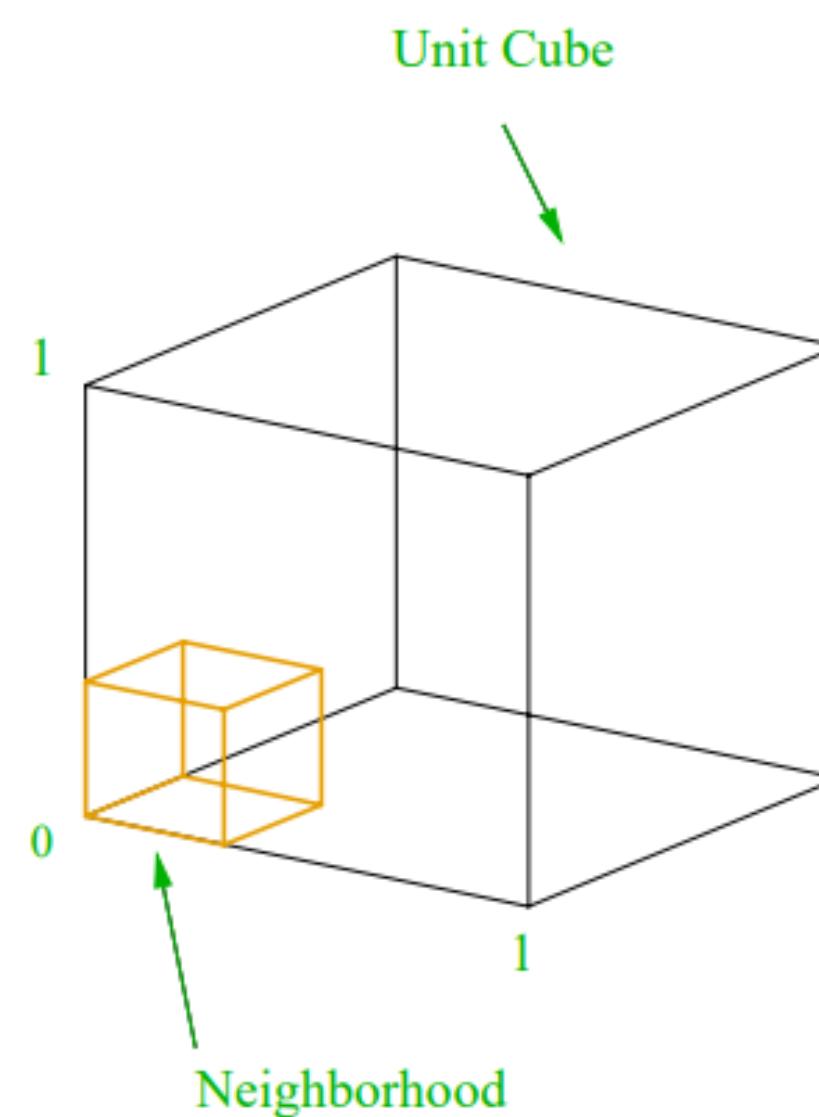
$$\begin{aligned}
 \hat{f}^{-\lambda}(t_n) &= b(t_n)^T \left[ (B^T W(t_0) B)^{-1} + \frac{k_\lambda(t_0, t_n)}{(-\gamma_{nn})} \cdot (B^T W(t_0) B)^{-1} b(t_n) b(t_n)^T (B^T W(t_0) B)^{-1} \right] \left[ B^T W(t_0) y - k_\lambda(t_0, t_n) b(t_n) y_n \right] \\
 &= \hat{f}(t_n) - b(t_n)^T (B^T W(t_0) B)^{-1} k_\lambda(t_0, t_n) b(t_n) y_n + \frac{\gamma_{nn} b(t_n)^T (B^T W(t_0) B)^{-1}}{(-\gamma_{nn})} \cdot B^T W(t_0) y - \frac{\gamma_{nn} b(t_n)^T (B^T W(t_0) B)^{-1} k_\lambda(t_0, t_n) b(t_n)}{(-\gamma_{nn})} y_n \\
 &= \hat{f}(t_n) - \gamma_{nn} y_n + \frac{\gamma_{nn} \hat{f}(t_n)}{(-\gamma_{nn})} - \frac{\gamma_{nn}^2 y_n}{(-\gamma_{nn})} \\
 &= \frac{\hat{f}(t_n) - \gamma_{nn} y_n}{(-\gamma_{nn})} \\
 \Rightarrow y_n - \hat{f}^{-\lambda}(t_n) &= y_n - \frac{\hat{f}(t_n) - \gamma_{nn} y_n}{(-\gamma_{nn})} = \frac{y_n - \hat{f}(t_n)}{(-\gamma_{nn})}.
 \end{aligned}$$

# Local Regression in $\mathbb{R}^p$

- Define the vector of polynomial terms in  $X$  of maximum degree  $d$ . For  $d = 1, p = 2$ ,  
 $b(X) = (1, X_1, X_2, X_1^2, X_2^2, X_1 X_2)$
- At each  $x_0 \in \mathbb{R}^p$ , solve  $\min_{\beta(x_0)} \sum_{i=1}^N K_\lambda(x_0, x_i) [y_i - b(x_i)^T \beta(x_0)]^2$  to produce the fit  
 $\hat{f}(x_0) = b(x_0)^T \hat{\beta}(x_0)$  with kernel metric  $\frac{\|x - x_0\|}{\lambda}$  where  $\|\cdot\|$  is the Euclidean norm.
- Since Euclidean norm depends on the units in each coordinate, it makes sense to standardize each predictor.

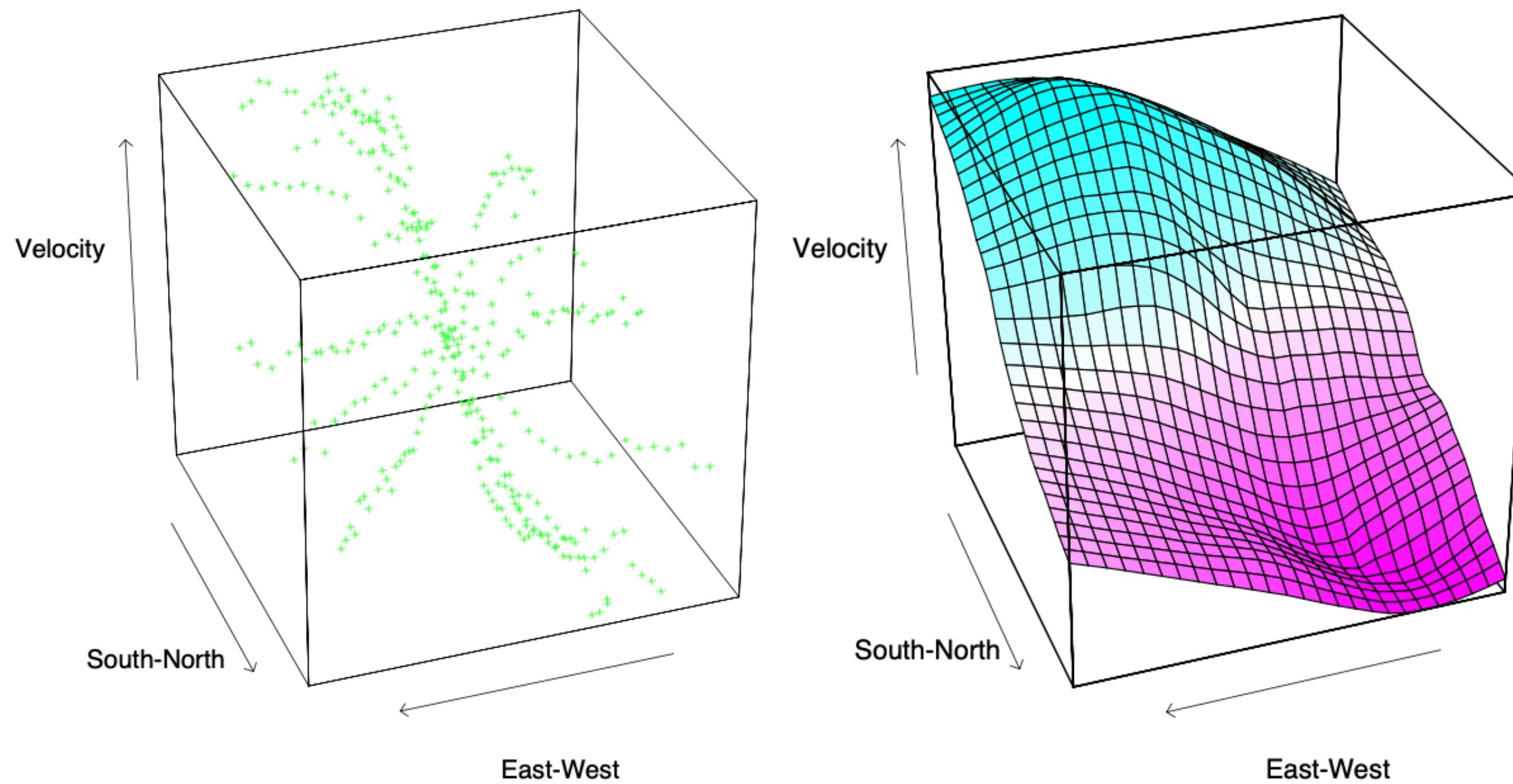
# Local Regression in $\mathbb{R}^p$

- **(Curse of dimensionality)** Boundary effects are a much higher problem in high dimensionality. Local polynomial regression seamlessly performs boundary correction to the desired order in any dimensions.



# Local Regression in $\mathbb{R}^p$

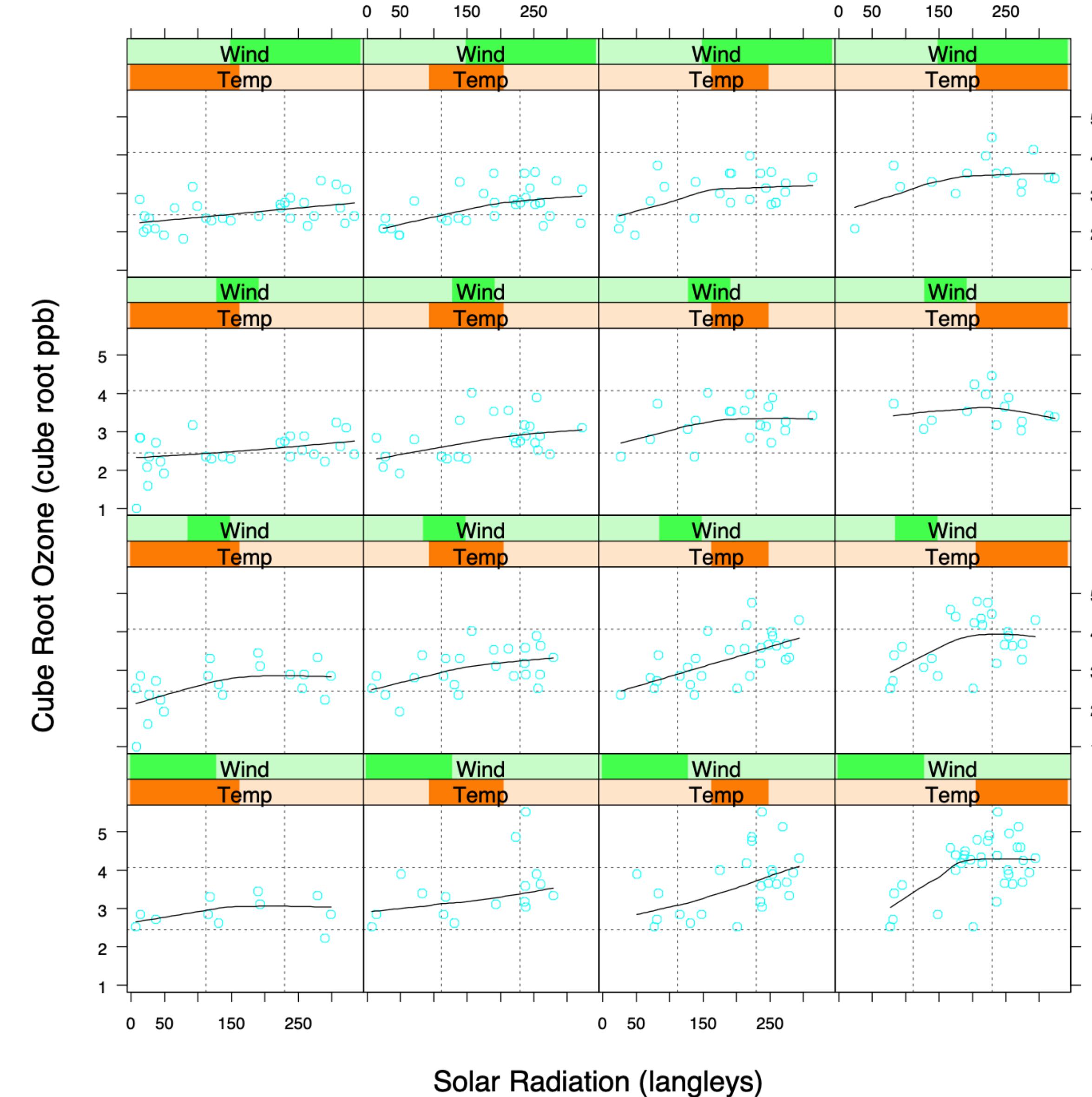
- **(Irregular Boundary case)** For  $\mathbb{R}^2$ , where the response is the velocity measurements on a galaxy, and the two predictors record positions on the celestial sphere.
- Local linear regression smoothing with NN window with 15% of the data



# Local Regression in $\mathbb{R}^p$

- **Structured Local Regression:** We need to make some structured assumptions about the model to avoid curse of dimensionality.
- Ex. For  $\mathbb{R}^3$  smoothing case. The response is ozone concentration with predictors temperature, wind speed, radiation. The trellis display shows ozone as a function of radiation, conditioned on intervals of temperature and wind speed.

# Local Regression in $\mathbb{R}^p$



# Local Regression in $\mathbb{R}^p$

- **ANOVA decompositions of the form:**  $\mathbb{E}(Y|X) = f(X_1, \dots, X_p) = \alpha + \sum_j g_j(X_j) + \sum_{k < l} g_{kl}(X_k, X_l) + \dots$
- **Additive models** assume only main effect terms:  $f(X) = \alpha + \sum_{j=1}^p g_j(X_j)$
- **Second order models:**  $f(X) = \alpha + \sum_j g_j(X_j) + \sum_{k < l} g_{kl}(X_k, X_l)$
- **Varying coefficients models:**  $f(X) = \alpha(Z) + \sum_{j=1}^q \beta_j(Z)X_j$  where  $Z = (X_{q+1}, \dots, X_p)$  with  $q < p$  and fitting by locally weighted least squares.

# Local Regression in $\mathbb{R}^p$

- **(Iterative back-fitting algorithms)** In the additive model, for example, if all but the  $k$ -th term is assumed known, then we can estimate  $\hat{g}_k$  based on

$$\{(x_i, y_i - \hat{\alpha} - \sum_{k \neq j} \hat{g}_j(X_j)) : i = 1, \dots, N\} . \text{ (Initialize } \hat{\alpha} = \bar{y}, \hat{g}_j(X_j) = 0 \text{) until}$$

convergence for  $j = 1, \dots, p$ . Notice that at any stage, one-dimensional local regression is all that is needed.

# Local Likelihood and Other Models

- Any parametric model can be made local if the fitting method accommodates observation weights.

- **Local Log Likelihood:** 
$$l(\beta(x_0)) = \sum_{i=1}^N K_\lambda(x_0, x_i) l(y_i, x_i^T \beta(x_0))$$

- **Ex. Locally weighted Logistic Regression:** 
$$l(\beta(x_0)) = \sum_{i=1}^N K_\lambda(x_0, x_i) \log p_{g_i}(x_i - x_0; \beta(x_0))$$

**(Centered at  $x_0$ )** where  $p_{g_i}(x_i; \beta(x_0)) = Pr(G = g_i | X = x)$ ,  $\log \frac{Pr(G = j | X = x)}{Pr(G = J | X = x)} = \beta_j(x_0)^T x$  for

$$j = 1, \dots, J - 1$$

# Local Likelihood and Other Models

Note that  $Pr(G = J | X = x) = \frac{1}{1 + \sum_{k=1}^{J-1} \exp(\beta_{k0}(x_0) + \beta_k(x_0)^T x)}$ ,

$$Pr(G = j | X = x) = \frac{\exp(\beta_{j0}(x_0) + \beta_j^T(x_0)x)}{1 + \sum_{k=1}^{J-1} \exp(\beta_{k0}(x_0) + \beta_k(x_0)^T x)} = Pr(G = J | X = x) \exp(\beta_{j0}(x_0) + \beta_j(x_0)^T x) \text{ for } j = 1, \dots, J-1$$

$$\begin{aligned} l(\beta(x_0)) &= \sum_{i=1}^N K_\lambda(x_0, x_i) \log p_{g_i}(x_i - x_0; \beta(x_0)) \\ &= \sum_{i=1}^N K_\lambda(x_0, x_i) \left[ \sum_{j=1}^{J-1} \left\{ I(y_i = j)(\beta_{j0}(x_0) + \beta_j(x_0)^T(x_i - x_0) + \log p_J(x_i - x_0; \beta_j(x_0))) \right\} + I(y_i = J) \log p_J(x_i - x_0; \beta_j(x_0)) \right] \\ &= \sum_{i=1}^N K_\lambda(x_0, x_i) \left[ \beta_{g_i 0} + \beta_{g_i}(x_0)^T(x_i - x_0) - \log \left\{ 1 + \sum_{k=1}^{J-1} \exp(\beta_{k0}(x_0) + \beta_k(x_0)^T(x_i - x_0)) \right\} \right] \end{aligned}$$

# Local Likelihood and Other Models

- We have centers the local regression at  $x_0$ , so that the fitted posterior probabilities at  $x_0$ :

$$\hat{Pr}(G = j | X = x_0) = \frac{\exp(\hat{\beta}_{j0}(x_0))}{1 + \sum_{k=1}^{J-1} \exp(\hat{\beta}_{k0}(x_0))}$$

- **Newton-Rapson Method:**

- Let  $X$  be centered at  $x_0$ ,  $W(x_0) = \text{diag}(K_\lambda(x_0, x_1), \dots, K_\lambda(x_0, x_N))$ ,  
 $W' = \text{diag}(K_\lambda(x_0, x_1)p(x_1 - x_0; \beta(x_0))(1 - p(x_1 - x_0; \beta(x_0))), \dots$   
 $\quad, K_\lambda(x_N, x_1)p(x_N - x_0; \beta(x_0))(1 - p(x_N - x_0; \beta(x_0)))^\top$

- Then  $\frac{\partial l(\beta(x_0))}{\partial \beta} = W(x_0)X^T(y - p)$ ,  $\frac{\partial l(\beta(x_0))}{\partial \beta \partial \beta^T} = -X^T W' X$

# Local Likelihood and Other Models

- $\beta(x_0)^{new} = \beta(x_0)^{old} + (X^T W' X)^{-1} W(x_0) X^T (y - p)$
- Since  $y_i \sim \text{Bernoulli}(p_i)$ , by CLT,  $\hat{\beta}(x_0) \sim \mathcal{N}(\beta(x_0), V)$  where  
 $V = (X^T W' X)^{-1} W(x_0) X^T W X W(x_0) (X^T W' X)^{-1}$ ,  
 $W = \text{diag}(p(x_1 - x_0; \beta(x_0))(1 - p(x_1 - x_0; \beta(x_0))), \dots,$   
 $p(x_N - x_0; \beta(x_0))(1 - p(x_N - x_0; \beta(x_0)))$
- Ex. (South Africa Heart Disease) Response variable is chd which is binary indicator, response is sbp and obesity respectively.
- $\text{logit} [Pr(\text{chd} | X)] = \beta_0(x_0)$ , Thus  $100(1 - \alpha)\%$  CI of  $f(x_0) : \hat{\beta}_0(x_0) \pm 2\sqrt{V_{11}}$

# Local Likelihood and Other Models

