

# Nonparametric Methods: I

Wei Li

Syracuse University

Spring 2024

# OVERVIEW

Nonparametric regression

Kernel smoothing (Locally weighted averages)

Local linear and polynomial regression

# Nonparametric regression

# Nonparametric regression

Given a random pair  $(X, Y) \in \mathbb{R}^p \times \mathbb{R}$ , the the **regression function** (of  $Y$  on  $X$  ) or **conditional expectation function**:

$$f(x) = E(Y \mid X = x)$$

► The basic goal is to estimate  $f$  from some i.i.d. sample  $(X_i, Y_i)_{i=1}^n$ .

For an i.i.d. sample  $(X_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}, i = 1, \dots, n,$ ,

$$Y_i = f(X_i) + \epsilon_i, \quad i = 1, \dots, n$$

-  $\epsilon_i, i = 1, \dots, n$  are i.i.d. random errors, with mean zero and satisfy  $E(\epsilon_i | X_i) = 0$ .

► nonparametric model does not assume  $f(x)$  take a known functional form

## example: k-nearest-neighbors regression

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

where  $N_k(x)$  contains the indices of the  $k$  closest points of  $x_1, \dots, x_n$  to  $x$ .

The regression prediction is a linear smoother (in  $\mathbf{y}$ ):

$$\hat{f}(x) = \sum_{i=1}^n w_i(x) y_i = \sum_{i=1}^n w(x, x_i) \cdot y_i = \mathbf{w}(x) \mathbf{y}$$

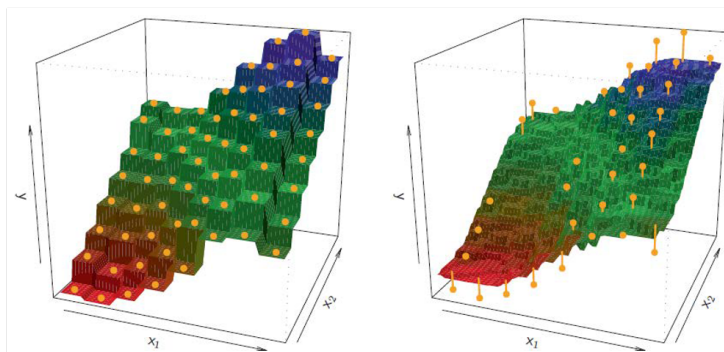
where the weights  $w_i(x), i = 1, \dots, n$  are defined as

$$w_i(x) = w(x, x_i) = 1(x_i \in N_k(x))/k,$$

Note: KNN is a memory-based model (lazy learner or exemplar-based learner).

# Curse of dimensionality

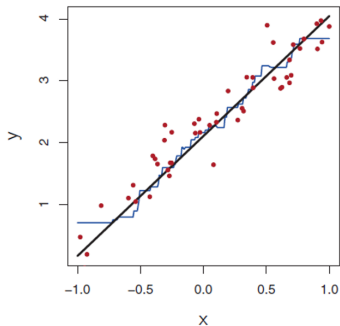
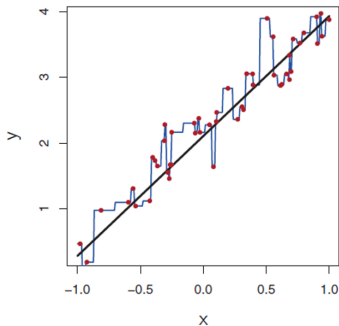
Suppose  $x \in \mathbb{R}^p$ .



$p=2$ .

- ▶ Left:  $k=1$ , most flexible fit, high variation
- ▶ Right:  $k=9$ , smoother fit, less variation

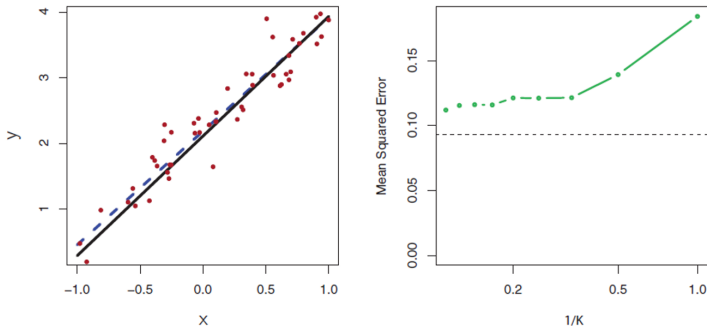
Source of Image: ISL Chapter 3.



$p=1$ , truth=linear

- ▶ left:  $k=1$
- ▶ right:  $k=9$

Source of Image: ISL Chapter 3.

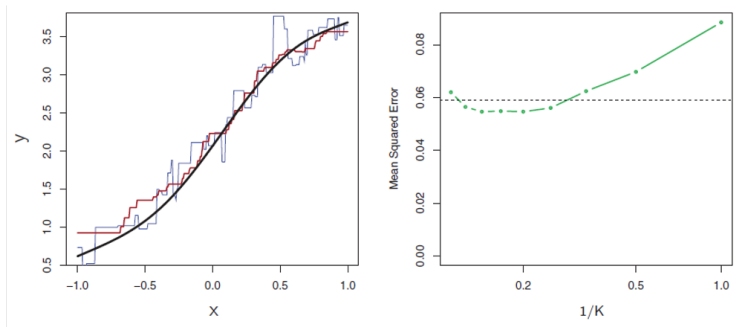


$p=1$ , truth=linear.

- ▶ Left: black(truth), dashed(LS fit)
- ▶ Right: dashed black (LS fit test MSE), dashed green (KNN test MSE)

Source of Image: ISL Chapter 3.

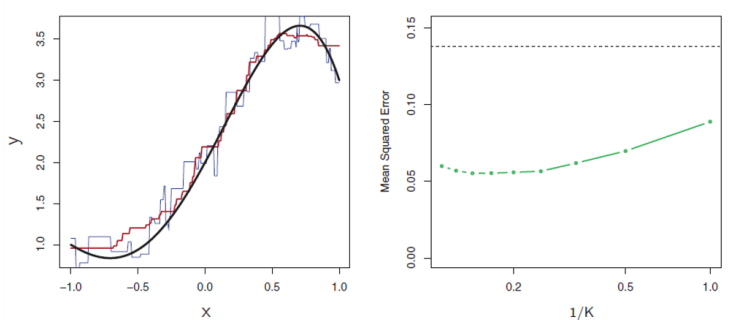




$p=1$ , truth=nonlinear.

- ▶ Left: black(truth), blue( $k=1$ ), red( $k=9$ )
- ▶ Right: dashed black (LS fit test MSE), dashed green (KNN test MSE)

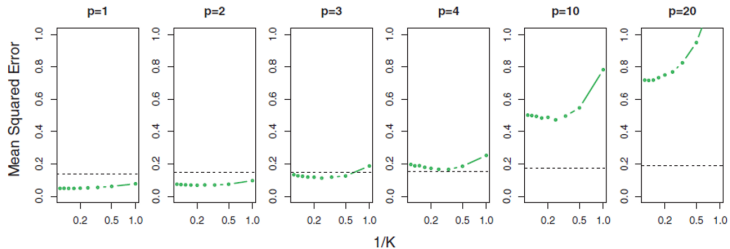
Source of Image: ISL Chapter 3.



$p=1$ , truth=nonlinear.

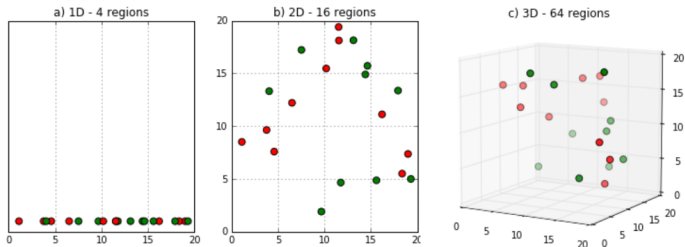
- ▶ Left: black(truth), blue( $k=1$ ), red( $k=9$ )
- ▶ Right: dashed black (LS fit test MSE), dashed green (KNN test MSE)

Source of Image: ISL Chapter 3.



If  $p$  is greater than 4, linear regression is superior to KNN even for nonlinear truth (sample size fixed)

Source of Image: ISL Chapter 3.



## Curse of Dimensionality

The **curse of dimensionality** refers to the principle that nonparametric estimation may become exponentially harder as the number of dimensions increases (with the exception of additional structure).

Source of Image:

<https://deepai.org/machine-learning-glossary-and-terms/curse-of-dimensionality>

## Convergence rates (k-NN regression)

$$\begin{aligned}\text{MSE}(\mathbf{x}_0) &= \mathbb{E}_{\tau} \left[ \hat{f}(\mathbf{x}_0) - f(\mathbf{x}_0) \right]^2 \\ &= \text{Var} \left( \frac{1}{k} \sum_{i \in N_k(\mathbf{x}_0)} y_i \right) + \left[ f(\mathbf{x}_0) - E \left( \frac{1}{k} \sum_{i \in N_k(\mathbf{x}_0)} y_i \right) \right]^2\end{aligned}$$

Assume  $f$  Lipschitz continuous

- ▶ variance term  $\approx \sigma^2/k$
- ▶ squared bias term  $\approx (k/n)^{2/p}$

With  $k \asymp n^{2/(2+p)}$ , it satisfies  $\text{MSE}(\mathbf{x}_0) \lesssim n^{-2/(2+p)}$

- ▶ optimal rate for estimating a Lipschitz  $p$ -variate functions
- ▶ requires exponentially more data to achieve a given error

# Kernel smoothing (Locally weighted averages)

## Kernel smoothing (Locally weighted averages)

Given a **bandwidth**  $h > 0$ , the **Nadaraya-Watson** kernel regression estimate is defined as

$$\hat{f}(x_0) = \frac{\sum_{i=1}^n K_h(x_0, x_i) y_i}{\sum_{i=1}^n K_h(x_0, x_i)}$$

Kernel smoothing is a linear smoothing method, with weights

$$w(x_0, x_i) = K_h(x_0, x_i) / \sum_{j=1}^n K_h(x_0, x_j).$$

That is  $\hat{f}(x_0) = \sum_{i=1}^n w(x_0, x_i) y_i$ .

$$\hat{\mathbf{f}} = (\hat{f}(x_1), \dots, \hat{f}(x_n))^{\top} = S\mathbf{y}$$

where  $S$  is  $n \times n$  matrix and  $S_{i,j} = w(x_i, x_j)$ .

## In general

$$K_h(x_0, x) = \bar{K}\left(\frac{\|x - x_0\|}{h}\right)$$

for some kernel function  $\bar{K} : \mathbb{R} \rightarrow \mathbb{R}$  satisfying

$$\int \bar{K}(t) dt = 1, \quad \int t \bar{K}(t) dt = 0, \quad 0 < \int t^2 \bar{K}(t) dt < \infty$$

The box-car kernel:

$$\bar{K}(t) = \begin{cases} 1 & |t| \leq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

the Gaussian kernel:

$$\bar{K}(t) = \frac{1}{\sqrt{2\pi}} \exp(-t^2/2)$$

the Epanechnikov kernel:

$$\bar{K}(t) = \begin{cases} 3/4 (1 - t^2) & \text{if } |t| \leq 1 \\ 0 & \text{else} \end{cases}$$

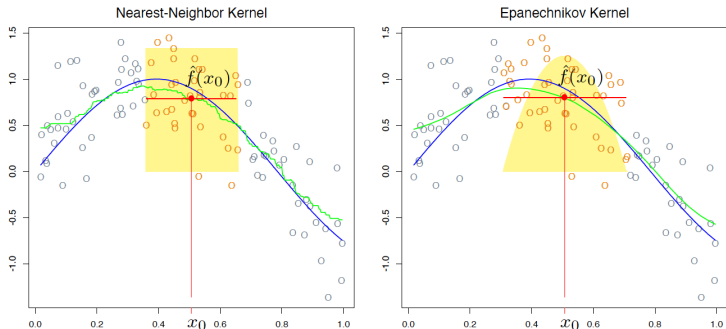


The k-nearest-neighbor estimator is just a raw (discontinuous) moving average of neighboring response with

$$K_h(x_0, x) = \bar{K}\left(\frac{\|x - x_0\|}{\|x_{(k)} - x_0\|}\right),$$

where

$$\bar{K}(t) = \begin{cases} 1 & |t| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$



The nearest-neighbor kernel compared with Epanechnikov kernel.

- ▶ blue curve is truth
- ▶ the red point is  $\hat{f}(x_0)$
- ▶ red circles are those observations contributing to the fit at  $x_0$
- ▶ solid yellow region indicates the weights assigned to observations.

ESL Figure 6.1

# Convergence rates (kernel regression)

In theory, for  $X \in \mathbb{R}^p$  and  $f$  is  $s$ -times differentiable ( $s \geq 1$ ), the minimal achievable rate of MSE by the best nonparametric method is  $n^{-2s/(2s+p)}$ .

- ▶ Suppose  $f$  is at least two-times differentiable ( $s = 2$ ), and  $X$  has a non-zero, differentiable density and the support is unbounded, using the Gaussian kernel (second order kernel), one can show

$$\begin{aligned}\text{Bias}(\hat{f}(x)) &\leq \tilde{C}_1 h^2 \\ \text{Var}(\hat{f}(x)) &\leq \frac{\tilde{C}_2}{nh^p}\end{aligned}$$

As long as  $h \rightarrow 0$  and  $nh^p \rightarrow \infty$ , the MSE goes to 0.

Balancing squared bias error and variance, the optimal choice is  $h = n^{-1/(p+4)}$  and the MSE is  $n^{-4/(p+4)}$ .

- ▶ the rate is “saturated” (bias rate of order not improve with  $s \geq 2$ ).
- ▶ poor bias rate near the boundaries if estimated on a bounded support

Is it possible to use a kernel estimate to obtain the optimal rate when  $s > 2$ ? The answer is yes, with the higher order kernel function  $\bar{K}$ .

A kernel function  $\bar{K}$  is said to be of order  $k$  provided that

$$\int \bar{K}(t)dt = 1, \int t^j \bar{K}(t)dt = 0, j = 1, \dots, k-1, \text{ and } 0 < \int t^k \bar{K}(t)dt < \infty$$

- ▶ Previous kernels were of order 2.
- ▶ An example of a 4 th-order kernel:  $\bar{K}(t) = \frac{3}{8} (3 - 5t^2) 1\{|t| \leq 1\}$ .
  - ▶ it can take negative values (a ubiquitous feature of higher-order kernel)

The second issue with the Gaussian kernel is related to the order  $O(h^2)$  of bias (assuming unbounded support)

- ▶ When estimated at the boundary points of the data, the bias becomes  $O(h)$ . This is the so-called “boundary effects”.
- ▶ One possible solution is to use some **boundary-corrected kernel**.

Perhaps better use the local polynomial regression, or series-function-based regression.

## General rates results

For kernel-based estimator, when  $f$  is  $s$ -times continuously differentiable ( $H^s$ -smoothness) and the higher-order kernel function is used, then

$$\sup_x |\hat{f}(x) - f(x)| = O(h^s) + O_p \left( \sqrt{\frac{\log n}{nh^p}} \right)$$

Balancing the two parts, the optimal choice  $h \asymp (\log n/n)^{1/(2s+p)}$ , yielding the sup-norm optimal rate  $(\log n/n)^{s/(2s+p)}$ .

- ▶ When  $h < (\log n/n)^{1/(2s+p)}$ , we call it **under-smoothing** (relative to the optimal rate)–bias dominant.
- ▶ When  $h > (\log n/n)^{1/(2s+p)}$ , we call it **over-smoothing** (relative to the optimal rate)–variance dominant.

under-smoothing/over-smoothing = underfitting/overfitting in asymptotic sense

# Bandwidth selection

1. rule of thumb:  $\hat{h} = \hat{c} \times (1/n)^{1/(4+p)}$  i.e.,  $s = 2$ , with estimated const.
2. cross-validation:

Choose from a sequence say,  $\{h_1, \dots, h_m\}$ ; randomly split the data into a  $K$ -parts partition  $(C_1, \dots, C_K)$

- ▶ For each part  $k$ :
  - ▶ For each value of  $h$ , fit the regression using the remaining  $K - 1$  parts (excluding  $k$ -th part), and denote the solution by  $\hat{f}_h^{-(k)}$ .
  - ▶ For each value of  $h$ , compute the prediction error on the  $k$ -th part,

$$CV_k(h) = \frac{1}{n_k} \sum_{i \in C_k} \left( y_i - \hat{f}_h^{-(k)}(x_i) \right)^2$$

- ▶ Take the average  $CV(h) = \frac{1}{K} \sum_{k=1}^K CV_k(h)$
- ▶ Find the best parameter  $\hat{h}$  which minimizes  $CV(h)$ .

One may repeat the whole process for a couple of times, obtaining  $CV(h)$ 's for different partitions, then take another average as the final CV estimate; and select  $h$  based on this final estimate.

# Inference about kernel estimation

For the kernel density estimator  $\hat{f}$

$$\hat{f}(x) - f(x) = (E\hat{f}(x) - f(x)) + (\hat{f}(x) - E\hat{f}(x))$$

- ▶  $B_n(x) = E\hat{f}(x) - f(x) = O(h^2)$  is the bias
- ▶  $\xi_n(x) = \hat{f}(x) - E\hat{f}(x) = O_p\left(\sqrt{1/nh^p}\right)$  is the stochastic variation and  $\sqrt{nh^p}\xi_n(x) \rightarrow N(0, V_n(x))$  for some variance  $V_n(x)$ .

A valid pointwise confidence interval for  $f(x)$  is

$$\left[ \hat{f}(x) - B_n(x) - z_{1-\alpha/2} \sqrt{\frac{V_n(x)}{nh^p}}, \hat{f}(x) - B_n(x) + z_{1-\alpha/2} \sqrt{\frac{V_n(x)}{nh^p}} \right]$$

where  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of a standard normal distribution.

# Inference about kernel estimation

$$\left[ \hat{f}(x) - B_n(x) - z_{1-\alpha/2} \sqrt{\frac{V_n(x)}{nh^p}}, \hat{f}(x) - B_n(x) + z_{1-\alpha/2} \sqrt{\frac{V_n(x)}{nh^p}} \right]$$

To eliminate bias  $B_n$ :

- ▶ *undersmooth* the kernel estimate, say  $h = o\left((nh^p)^{-1/4}\right)$ , or
- ▶ estimate  $B_n$

Without undersmoothing, ignoring  $B_n$  will in general produce a biased pointwise confidence interval

\*Since the kernel regression is also linear smoother, can use the formula for the pointwise CI for series estimation.



# Confidence bands

It is possible to obtain (sup-norm) confidence band (say through bootstrap). Suppose  $R_\alpha$  is the  $1 - \alpha$  quantile of  $\sup_x |\hat{f}(x) - f(x)|$ , then

$$P\left(f(x) \in (\hat{f}(x) - R_\alpha, \hat{f}(x) + R_\alpha), \forall x\right) = 1 - \alpha.$$

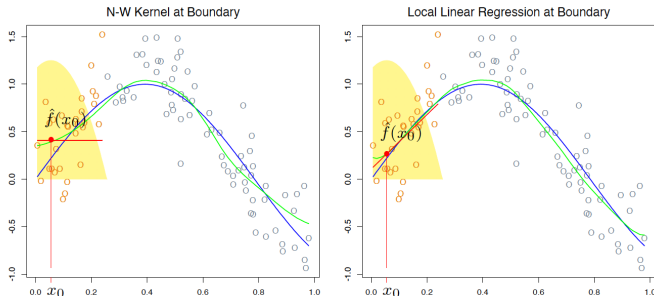
So a  $1 - \alpha$  uniform confidence band for  $f(x)$  is  $(\hat{f}(x) - R_\alpha, \hat{f}(x) + R_\alpha)$ .

One can use bootstrap for  $\sup_x |\hat{f}(x) - E(\hat{f}(x))|$  to address the variance estimate, using debiased estimator or undersmoothing to address the bias term. See future lessons.

## Local linear and polynomial regression

# Local linear and polynomial regression

The kernel regression suffers from poor bias at the boundaries of the domain of the inputs  $x_1, \dots, x_n$ .



NW estimator compared with local linear regression.

ESL Figure 6.1

## Locally linear (locally weighted) regression: $p = 1$

At each point, we make prediction by using linear regression weighting only nearby points. At each target point  $x_0$ , solve

$$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^n K_h(x_0, x_i) [y_i - \alpha(x_0) - \beta(x_0) x_i]^2$$

- ▶ vector-valued function  $b(x)^\top = (1, x)$ .
- ▶  $\mathbf{B}$  be the  $n \times 2$  regression matrix with  $i$ th row  $b(x_i)^\top$ ,
- ▶  $\mathbf{W}(x_0)$  the  $n \times n$  diagonal matrix with  $i$ th diagonal element  $K_h(x_0, x_i)$

$$\begin{aligned}\hat{f}(x_0) &= b(x_0)^\top (\mathbf{B}^\top \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{W}(x_0) \mathbf{y} \\ &= \sum_{i=1}^n l_i(x_0) y_i = l(x_0)^\top \mathbf{y}\end{aligned}$$

A local polynomial of any degree  $d$ :

$$\min_{\alpha(x_0), \beta_j(x_0), j=1, \dots, d} \sum_{i=1}^n K_h(x_0, x_i) \left[ y_i - \alpha(x_0) - \sum_{j=1}^d \beta_j(x_0) x_i^j \right]^2$$

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

$$\hat{f}(x_0) = b(x_0) (\mathbf{B}^\top \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{W}(x_0) \mathbf{y} = l(x_0)^\top \mathbf{y}$$

- ▶  $b(x) = (1, x, \dots, x^d)$
- ▶  $\mathbf{B}$  is an  $n \times (d+1)$  matrix with  $i$ -th row  $b(x_i) = (1, x_i, \dots, x_i^d)$ ,
- ▶  $\mathbf{W}(x_0)$  is as before.

Note local polynomial regression is a linear smoother, so  $\hat{f} = S_\lambda \mathbf{y}$  where  $[S_\lambda]_{i,j} = l_j(x_i)$ .

$p > 1$

Let  $b(X)$  be a vector of polynomial terms in  $X$  of maximum degree  $d$ .

- For example, with  $d = 1$  and  $p = 2$  we get  $b(X) = (1, X_1, X_2)$ ;  
with  $d = 2$  we get  $b(X) = (1, X_1, X_2, X_1^2, X_2^2, X_1X_2)$

At each  $x_0 \in \mathbb{R}^p$  solve

$$\min_{\beta(x_0)} \sum_{i=1}^n K_h(x_0, x_i) \left( y_i - b(x_i)^\top \beta(x_0) \right)^2$$

to produce the fit  $\hat{f}(x_0) = b(x_0)^\top \hat{\beta}(x_0)$ .

The polynomial terms are constructed as tensor products of basis functions. See series estimators.

# Local likelihood

Associated with each observation  $y_i$  is a parameter  $\theta_i = \theta(x_i) = x_i^\top \beta$  linear in the covariate(s)  $x_i$ , and inference for  $\beta$  is based on the loglikelihood  $\ell(\beta) = \sum_{i=1}^n \log L(y_i, x_i^\top \beta)$ .

Model  $\theta(x)$  by using the likelihood local to  $x_0$  for inference of  $\theta(x_0) = x_0^\top \beta(x_0)$  :

$$\ell(\beta(x_0)) = \sum_{i=1}^n K_h(x_0, x_i) \log L(y_i, x_i^\top \beta(x_0))$$

# Local version of multi-class linear logistic regression model

The data consist of features  $x_i$  and an associated categorical response  $y_i \in \{1, 2, \dots, J\}$ , and the linear model has the form

$$\Pr(Y = j \mid X = x) = \frac{e^{\beta_{j0} + \beta_j^\top x}}{1 + \sum_{k=1}^{J-1} e^{\beta_{k0} + \beta_k^\top x}}$$

where  $\beta_{J0} := 0$  and  $\beta_J := 0$  by the definition of the model.

The local log-likelihood for this  $J$  class model can be written

$$\sum_{i=1}^n K_h(x_0, x_i) \left\{ \beta_{y_i 0}(x_0) + \beta_{y_i}(x_0)^\top (x_i - x_0) - \log \left[ 1 + \sum_{k=1}^{J-1} \exp \left( \beta_{k0}(x_0) + \beta_k(x_0)^\top (x_i - x_0) \right) \right] \right\}$$