
Supplementary material for “Robust modal regression with direct gradient approximation of modal regression risk”

Hiroaki Sasaki¹, Tomoya Sakai², Takafumi Kanamori^{3,4}

¹Future University Hakodate, Hokkaido, Japan ²NEC Corporation, Tokyo, Japan

³Tokyo Institute of Technology, Tokyo, Japan ⁴RIKEN AIP, Tokyo, Japan

A Proof of Theorem 1

Proof. Let us denote the inner product in a reproducing kernel Hilbert space \mathcal{H} (RKHS) by $\langle \cdot, \cdot \rangle_{\mathcal{H}}$. Since the empirical Fisher divergence (11) can be expressed as

$$\widehat{J}(r) = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \langle r, k(\cdot, \mathbf{z}_i) \rangle_{\mathcal{H}}^2 + \langle r, \partial'_y k(\cdot, \mathbf{z}_i) \rangle_{\mathcal{H}} \right], \quad (28)$$

where k is the kernel function in \mathcal{H} , the representer theorem for derivatives [Zhou, 2008] ensures that the model r should take the following optimal form:

$$r(\mathbf{z}) = \sum_{i=1}^n [\alpha_i k(\mathbf{z}, \mathbf{z}_i) + \beta_i \partial'_y k(\mathbf{z}, \mathbf{z}_i)], \quad (29)$$

where $\mathbf{z} = (y, \mathbf{x}^\top)^\top$, $\partial'_y k(\mathbf{z}, \mathbf{z}_i) := \frac{\partial}{\partial y'} k(\mathbf{z}, \mathbf{z}')|_{\mathbf{z}'=\mathbf{z}_i}$ with $\mathbf{z}' = (y', \mathbf{x}'^\top)^\top$ (i.e., ∂'_y denotes the partial derivative with respect to the second variable of the kernel function k), and α_i and β_i are coefficients to be estimated. Computing the partial derivative of (29) with respect to y yields

$$\partial_y \widehat{r}(\mathbf{z}) = \sum_{i=1}^n [\alpha_i \partial_y k(\mathbf{z}, \mathbf{z}_i) + \beta_i \partial_y \partial'_y k(\mathbf{z}, \mathbf{z}_i)], \quad (30)$$

where $\partial_y := \frac{\partial}{\partial y}$.

Next, we define the (i, j) -th element in matrices, \mathbf{K} , \mathbf{G} and \mathbf{H} , by

$$[\mathbf{K}]_{ij} = k(\mathbf{z}_i, \mathbf{z}_j), \quad [\mathbf{G}]_{ij} = \partial'_y k(\mathbf{z}_i, \mathbf{z}_j) \quad \text{and} \quad [\mathbf{H}]_{ij} = \partial_y \partial'_y k(\mathbf{z}_i, \mathbf{z}_j).$$

Then, $\mathbf{r} = (r(\mathbf{z}_1), r(\mathbf{z}_2), \dots, r(\mathbf{z}_n))^\top$ and $\partial_y \mathbf{r} = (\partial_y r(\mathbf{z}_1), \partial_y r(\mathbf{z}_2), \dots, \partial_y r(\mathbf{z}_n))^\top$ are compactly expressed as

$$\mathbf{r} = \mathbf{K} \boldsymbol{\alpha} + \mathbf{G} \boldsymbol{\beta} \quad (31)$$

$$\partial_y \mathbf{r} = \mathbf{G}^\top \boldsymbol{\alpha} + \mathbf{H} \boldsymbol{\beta}. \quad (32)$$

Regarding the RKHS norm,

$$\begin{aligned} \|r\|_{\mathcal{H}}^2 &= \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} + 2 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \beta_j \partial'_y k(\mathbf{z}_i, \mathbf{z}_j) + \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j \partial_y \partial'_y k(\mathbf{z}_i, \mathbf{z}_j) \\ &= \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} + 2 \boldsymbol{\alpha}^\top \mathbf{G} \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{H} \boldsymbol{\beta}, \end{aligned} \quad (33)$$

Substituting (31), (32) and (33) into (28) yields

$$\begin{aligned}\tilde{J}(r) &:= \hat{J}(r) + \frac{\lambda}{2} \|r\|_{\mathcal{H}}^2 \\ &= \frac{1}{2n} \|\mathbf{K}\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta}\|^2 + \frac{1}{n} \mathbf{1}_n^\top (\mathbf{G}^\top \boldsymbol{\alpha} + \mathbf{H}\boldsymbol{\beta}) + \frac{\lambda}{2} (\boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha} + 2\boldsymbol{\alpha}^\top \mathbf{G}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{H}\boldsymbol{\beta}).\end{aligned}$$

Taking the derivatives of \tilde{J} with respect to $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ yields

$$\begin{aligned}\frac{\partial \tilde{J}(r)}{\partial \boldsymbol{\alpha}} &= \frac{1}{n} \mathbf{K}(\mathbf{K}\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta}) + \frac{1}{n} \mathbf{G}\mathbf{1}_n + \lambda \mathbf{K}\boldsymbol{\alpha} + \lambda \mathbf{G}\boldsymbol{\beta} \\ &= \frac{1}{n} \mathbf{K} \{(\mathbf{K} + n\lambda \mathbf{I}_n)\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta}\} + \mathbf{G} \left\{ \frac{1}{n} \mathbf{1}_n + \lambda \boldsymbol{\beta} \right\} \\ \frac{\partial \tilde{J}(r)}{\partial \boldsymbol{\beta}} &= \frac{1}{n} \mathbf{G}^\top (\mathbf{K}\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta}) + \frac{1}{n} \mathbf{H}\mathbf{1}_n + \lambda \mathbf{H}\boldsymbol{\beta} + \lambda \mathbf{G}^\top \boldsymbol{\alpha} \\ &= \frac{1}{n} \mathbf{G}^\top \{(\mathbf{K} + n\lambda \mathbf{I}_n)\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta}\} + \mathbf{H} \left\{ \frac{1}{n} \mathbf{1}_n + \lambda \boldsymbol{\beta} \right\}.\end{aligned}$$

The optimality condition is given by

$$(\mathbf{K} + n\lambda \mathbf{I}_n)\boldsymbol{\alpha} + \mathbf{G}\boldsymbol{\beta} = \mathbf{0}, \quad \frac{1}{n} \mathbf{1}_n + \lambda \boldsymbol{\beta} = \mathbf{0}.$$

Thus, the optimal coefficients can be computed as

$$\hat{\boldsymbol{\alpha}} = \frac{1}{n\lambda} (\mathbf{K} + n\lambda \mathbf{I}_n)^{-1} \mathbf{G}\mathbf{1}_n, \quad \hat{\boldsymbol{\beta}} = -\frac{1}{n\lambda} \mathbf{1}_n.$$

Substituting $\hat{\boldsymbol{\alpha}}$ and $\hat{\boldsymbol{\beta}}$ into (29) completes the proof. \square

B Leave-one-out cross-validation

Here, we show that the LOOCV score can be efficiently computed by following Kanamori et al. [2012]. The notations in Section A are inherited

Let us denote the collection of data samples except $\mathbf{z}_l := (y_l, \mathbf{x}_l^\top)^\top$ by \mathcal{D}_l (i.e., $\mathcal{D} \setminus \mathbf{z}_l$). K-LSLD from \mathcal{D}_l is given by

$$\hat{r}^{(l)}(\mathbf{z}) := \sum_{\substack{i=1 \\ i \neq l}}^n \left[\hat{\alpha}_i^{(l)} k(\mathbf{z}, \mathbf{z}_i) + \hat{\beta}_i^{(l)} \partial'_y k(\mathbf{z}, \mathbf{z}_i) \right],$$

where

$$\hat{\boldsymbol{\alpha}}^{(l)} := \frac{1}{(n-1)\lambda} (\mathbf{K}^{(l)} + (n-1)\lambda \mathbf{I}_{n-1})^{-1} \mathbf{G}^{(l)} \mathbf{1}_{n-1}, \quad \hat{\boldsymbol{\beta}}^{(l)} := -\frac{1}{(n-1)\lambda} \mathbf{1}_{n-1}.$$

In the equations above, $\mathbf{K}^{(l)}$ and $\mathbf{G}^{(l)}$ are \mathbf{K} and \mathbf{G} except \mathbf{z}_l , respectively. Then, the LOOCV score can be computed as

$$\text{LOOCV} = \frac{1}{n} \sum_{l=1}^n \left[\frac{1}{2} \{\hat{r}^{(l)}(\mathbf{z}_l)\}^2 + \partial'_y \hat{r}^{(l)}(\mathbf{z}_l) \right].$$

However, to naively compute this LOOCV score, we need to compute the inverse of $n-1$ by $n-1$ matrix for all $\hat{\boldsymbol{\alpha}}^{(l)}$, $l = 1, \dots, n$, which is time-consuming.

To cope with this problem, we derive an equivalent form of $\hat{r}^{(l)}$. $\hat{\boldsymbol{\alpha}}^{(l)}$ can be regarded as the solution of the optimization problem,

$$\hat{\boldsymbol{\alpha}}^{(l)} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^{n-1}}{\text{argmin}} \left[\frac{1}{2} \boldsymbol{\alpha}^\top (\mathbf{K}^{(l)} + (n-1)\lambda \mathbf{I}_{n-1}) \boldsymbol{\alpha} - \frac{1}{(n-1)\lambda} \mathbf{1}_{n-1}^\top \mathbf{G}^{(l)} \boldsymbol{\alpha} \right].$$

Here, we solve an alternative optimization problem as

$$\tilde{\boldsymbol{\alpha}}^{(l)} := \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\operatorname{argmin}} \left[\frac{1}{2} \boldsymbol{\alpha}^\top (\mathbf{K} + (n-1)\lambda \mathbf{I}_n) \boldsymbol{\alpha} - \frac{1}{(n-1)\lambda} (\mathbf{1}_n - \mathbf{e}_l)^\top \mathbf{G} \boldsymbol{\alpha} \right] \quad \text{s.t.} \quad \tilde{\alpha}_l^{(l)} = 0, \quad (34)$$

where \mathbf{e}_l is the unit vector with the l -th element being 1. Note that $\tilde{\boldsymbol{\alpha}}^{(l)} \in \mathbb{R}^n$ is an n -dimensional vector, while $\hat{\boldsymbol{\alpha}}^{(l)} \in \mathbb{R}^{n-1}$ is an $(n-1)$ -dimensional one. With $\tilde{\boldsymbol{\alpha}}^{(l)} = (\tilde{\alpha}_1^{(l)}, \dots, \tilde{\alpha}_n^{(l)})^\top$, $\hat{\mathbf{r}}^{(l)}$ can be equivalently expressed as

$$\hat{\mathbf{r}}^{(l)}(\mathbf{z}) = \sum_{i=1}^n \left[\tilde{\alpha}_i^{(l)} k(\mathbf{z}, \mathbf{z}_i) + \tilde{\beta}_i^{(l)} \partial_y' k(\mathbf{z}, \mathbf{z}_i) \right], \quad (35)$$

where

$$\tilde{\boldsymbol{\beta}}^{(l)} = (\tilde{\beta}_1^{(l)}, \dots, \tilde{\beta}_n^{(l)})^\top := -\frac{1}{(n-1)\lambda} (\mathbf{1}_n - \mathbf{e}_l). \quad (36)$$

Applying the method of Lagrange multipliers to (34) yields the following solution:

$$\tilde{\boldsymbol{\alpha}}^{(l)} = (\mathbf{K} + (n-1)\lambda \mathbf{I}_n)^{-1} \left\{ \frac{1}{(n-1)\lambda} \mathbf{G} (\mathbf{1}_{n-1} - \mathbf{e}_l) + t_l \mathbf{e}_l \right\}, \quad (37)$$

where t_l is fixed such that $\tilde{\alpha}_l^{(l)} = 0$. Eq.(37) clearly shows that in order to obtain $\tilde{\boldsymbol{\alpha}}^{(l)}$ for all $l = 1, \dots, n$, it is sufficient to compute the inverse of the n by n matrix only once, which significantly reduces the cost to compute the LOOCV score compared with the naive way discussed above.

Next, we show a compact and analytic form of the LOOCV score. To this end, we define

$$\mathbf{A} := (\tilde{\boldsymbol{\alpha}}^{(1)}, \dots, \tilde{\boldsymbol{\alpha}}^{(n)}) \quad \text{and} \quad \mathbf{B} := (\tilde{\boldsymbol{\beta}}^{(1)}, \dots, \tilde{\boldsymbol{\beta}}^{(n)}).$$

Then, from (36) and (37), \mathbf{A} and \mathbf{B} can be computed as

$$\mathbf{A} = \mathbf{L}(\mathbf{S} - \mathbf{T}) \quad \text{and} \quad \mathbf{B} = -\frac{1}{(n-1)\lambda} \mathbf{E}, \quad (38)$$

where $\mathbf{L} := (\mathbf{K} + (n-1)\lambda \mathbf{I}_n)^{-1}$, $\mathbf{S} := \frac{1}{(n-1)\lambda} \mathbf{G} \mathbf{E}$,

$$[\mathbf{E}]_{ij} := \begin{cases} 0 & i = j, \\ 1 & i \neq j, \end{cases} \quad \text{and} \quad [\mathbf{T}]_{ij} := \begin{cases} [\mathbf{L}\mathbf{S}]_{ii}/[\mathbf{L}]_{ii} & i = j, \\ 0 & i \neq j. \end{cases}$$

Finally, the LOOCV score can be computed analytically as

$$\text{LOOCV} = \frac{1}{n} \left\{ \frac{1}{2} \tilde{\mathbf{r}}^\top \tilde{\mathbf{r}} + \mathbf{1}_n^\top \partial_y \tilde{\mathbf{r}} \right\}, \quad (39)$$

where

$$\begin{aligned} \tilde{\mathbf{r}} &:= (\hat{\mathbf{r}}^{(1)}(\mathbf{z}_1), \hat{\mathbf{r}}^{(2)}(\mathbf{z}_2), \dots, \hat{\mathbf{r}}^{(n)}(\mathbf{z}_n))^\top = (\mathbf{K} \odot \mathbf{A}^\top + \mathbf{G} \odot \mathbf{B}^\top) \mathbf{1}_n \\ \partial_y \tilde{\mathbf{r}} &:= (\partial_y \hat{\mathbf{r}}^{(1)}(\mathbf{z}_1), \partial_y \hat{\mathbf{r}}^{(2)}(\mathbf{z}_2), \dots, \partial_y \hat{\mathbf{r}}^{(n)}(\mathbf{z}_n))^\top = (\mathbf{G}^\top \odot \mathbf{A}^\top + \mathbf{H} \odot \mathbf{B}^\top) \mathbf{1}_n. \end{aligned}$$

The symbol \odot denotes element-wise multiplication.

C Proof of Theorem 2

With the assumption that $\hat{\alpha}_l = 0$ for all l , we express $\hat{D}_{\hat{\mathbf{r}}}[\boldsymbol{\theta}_2 | \boldsymbol{\theta}_1]$ as

$$\begin{aligned} \hat{D}_{\hat{\mathbf{r}}}[\boldsymbol{\theta}_2 | \boldsymbol{\theta}_1] &= \frac{1}{n} \sum_{i=1}^n \int_0^1 \hat{\mathbf{r}}(\boldsymbol{\theta}(t)^\top \mathbf{k}_m(\mathbf{x}_i), \mathbf{x}_i) \mathbf{k}_m(\mathbf{x}_i)^\top (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1) dt \\ &= \frac{1}{n} \sum_{i,l=1}^n \underbrace{\left[\int_0^1 \frac{y_l - \boldsymbol{\theta}(t)^\top \mathbf{k}_m(\mathbf{x}_i)}{n\lambda\sigma_y^2} \varphi \left\{ \frac{(\boldsymbol{\theta}(t)^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2} \right\} \mathbf{k}_m(\mathbf{x}_i)^\top (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1) dt \right]}_{(*)} k_x(\mathbf{x}_i, \mathbf{x}_l), \quad (40) \end{aligned}$$

By the substitution $Y_l = \frac{y_l - \boldsymbol{\theta}^\top(t) \mathbf{k}_m(\mathbf{x}_i)}{\sigma_y}$, the integral (\star) in (40) is computed as

$$(\star) = -\frac{1}{n\lambda} \int_{Y_l^{(1)}}^{Y_l^{(2)}} Y_l \varphi\left(\frac{Y_l^2}{2}\right) dY_l = \frac{1}{n\lambda} \left[\phi\left\{\frac{(\boldsymbol{\theta}_2^\top \mathbf{k}(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2}\right\} - \phi\left\{\frac{(\boldsymbol{\theta}_1^\top \mathbf{k}(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2}\right\} \right], \quad (41)$$

where we used $\frac{dY_l}{dt} = \frac{(\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1)^\top \mathbf{k}_m(\mathbf{x}_i)}{\sigma_y}$ from the path (24), $Y_l^{(1)} = \frac{y_l - \boldsymbol{\theta}_1^\top \mathbf{k}_m(\mathbf{x}_i)}{\sigma_y}$ and $Y_l^{(2)} = \frac{y_l - \boldsymbol{\theta}_2^\top \mathbf{k}_m(\mathbf{x}_i)}{\sigma_y}$.

Then, substituting (41) into (40) yields

$$\begin{aligned} \widehat{D}_{\widehat{r}}[\boldsymbol{\theta}_2 | \boldsymbol{\theta}_1] &= \frac{1}{n^2 \lambda} \sum_{i,l=1}^n k_x(\mathbf{x}_i, \mathbf{x}_l) \left[\phi\left\{\frac{(\boldsymbol{\theta}_2^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2}\right\} - \phi\left\{\frac{(\boldsymbol{\theta}_1^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2}\right\} \right] \\ &\geq \frac{1}{n^2 \lambda} \sum_{i,l=1}^n k_x(\mathbf{x}_i, \mathbf{x}_l) \varphi\left\{\frac{(\boldsymbol{\theta}_1^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2}\right\} \left\{ \frac{(\boldsymbol{\theta}_1^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2} - \frac{(\boldsymbol{\theta}_2^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2\sigma_y^2} \right\} \\ &= \frac{1}{2} \{ \boldsymbol{\theta}_1^\top \mathbf{H}(\boldsymbol{\theta}_1) \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2^\top \mathbf{H}(\boldsymbol{\theta}_1) \boldsymbol{\theta}_2 - 2(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^\top \mathbf{h}(\boldsymbol{\theta}_1) \}, \end{aligned}$$

where we applied a well-known inequality for convex functions as

$$\phi(t_2) - \phi(t_1) \geq \varphi(t_1)(t_1 - t_2),$$

where $\varphi(t) := -\frac{d}{dt}\phi(t)$.

By $\boldsymbol{\theta}_1 \leftarrow \boldsymbol{\theta}^\tau$ and $\boldsymbol{\theta}_2 \leftarrow \boldsymbol{\theta}^{\tau+1}$, we have

$$\begin{aligned} \widehat{D}_{\widehat{r}}[\boldsymbol{\theta}^{\tau+1} | \boldsymbol{\theta}^\tau] &\geq \frac{1}{2} \{ \boldsymbol{\theta}^{\tau+1 \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1 \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1} - 2(\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1})^\top \mathbf{h}(\boldsymbol{\theta}^\tau) \} \\ &= \frac{1}{2} \{ \boldsymbol{\theta}^{\tau \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1 \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1} - 2(\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1})^\top \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1} \} \\ &= \frac{1}{2} \{ \boldsymbol{\theta}^{\tau \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^\tau + \boldsymbol{\theta}^{\tau+1 \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1} - 2\boldsymbol{\theta}^{\tau \top} \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1} \} \\ &= \frac{1}{2} (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1})^\top \mathbf{H}(\boldsymbol{\theta}^\tau) (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau+1}), \end{aligned}$$

where we used the relation $\mathbf{h}(\boldsymbol{\theta}^\tau) = \mathbf{H}(\boldsymbol{\theta}^\tau) \boldsymbol{\theta}^{\tau+1}$ in (20) on the first line. Since $\mathbf{H}(\boldsymbol{\theta})$ is assumed to be positive definite, the proof is completed.

D Regression methods in Section 5.1

Section 5.1 estimates the conditional mode by the following kernel model:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{k}_m(\mathbf{x}) = \sum_{i=1}^n \theta_i k_m(\mathbf{x}, \mathbf{x}_i),$$

where $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^\top$ is the vector of parameters, $\mathbf{k}_m(\mathbf{x}) = (k_m(\mathbf{x}, \mathbf{x}_1), \dots, k_m(\mathbf{x}, \mathbf{x}_n))$ and $k_m(\mathbf{x}, \mathbf{y})$ is a kernel function. We employed the Gaussian kernel for $k_m(\mathbf{x}, \mathbf{y})$ where the width parameter was fixed at the median of the pairwise distance $\|\mathbf{x}_i - \mathbf{x}_j\|$ (i.e., the median trick) as done in Gretton et al. [2012]. The following regression methods were applied to the same artificial datasets:

- *Kernel ridge regression (KRR)*: $f_{\boldsymbol{\theta}}(\mathbf{x})$ was estimated under the squared-loss with the RKHS norm regularization as follows:

$$\min_{\boldsymbol{\theta}} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2 + \lambda \|f_{\boldsymbol{\theta}}\|_{\mathcal{H}}^2 \right].$$

The regularization parameter λ was determined by the five-fold cross-validation.

- *Least absolute deviations (LAD)*: The absolute deviation was used as the loss function:

$$\min_{\boldsymbol{\theta}} \left[\frac{1}{n} \sum_{i=1}^n |y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i)| + \lambda \|f_{\boldsymbol{\theta}}\|_{\mathcal{H}}^2 \right].$$

As in [Feng et al., 2017, Algorithm 1], the iteratively reweighted least squares (IRLS) algorithm was applied to optimize the parameters. The five-fold cross-validation was performed to select the regularization parameter.

- *Huber loss (Huber)*: $f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{k}_m(\mathbf{x})$ was estimated with the Huber loss by minimizing

$$\min_{\boldsymbol{\theta}} \left[\frac{1}{n} \sum_{i=1}^n L(y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \lambda \|f_{\boldsymbol{\theta}}\|_{\mathcal{H}}^2 \right],$$

where with a positive parameter γ ,

$$L(t) = \begin{cases} \frac{t^2}{2\gamma}, & |t| \leq \gamma, \\ |t| - \gamma/2, & |t| > \gamma, \end{cases}$$

Following [Feng et al., 2017, Algorithm 1], again, IRLS was used to optimize the model parameters. The five-fold cross-validation was performed to select γ and λ using the LAD loss¹.

- *Variational heteroscedastic Gaussian process regression (VHGPR)* [Lázaro-Gredilla and Titsias, 2011]: We used the MATLAB code, which is available at <http://www.tsc.uc3m.es/~miguel/downloads.php>.
- *Modal regression with kernel density estimation (MR_{KDE})*: A variant of DMR-K with kernel density estimation (KDE) following the naive two-step approach. As done in Yao et al. [2012], KDE was performed to estimate the joint density $p_{y\mathbf{x}}(y, \mathbf{x})$ where the Gauss kernel was employed and the width parameters in the kernel were determined by the standard least-squares cross-validation [Wasserman, 2006]. A similar update rule as DMR-K was derived and used similarly as in Algorithm 1. Details are given in Section E.
- *Direct modal regression with kernels (DMR-K)*: A proposed method based on reproducing kernels. Regarding K-LSLD, the Gaussian kernel was used both for k_x and k_y , and the width parameter in each kernel is determined by the leave-one-out cross-validation method in Section 3, while we fixed the regularization parameter at $n^{-0.9}$ by following Kanamori et al. [2012]. Then, $f_{\boldsymbol{\theta}}(\mathbf{x})$ was estimated according to Algorithm 1.

Regarding both MR_{KDE} and DMR-K, we initialized the parameters $\boldsymbol{\theta}$ by LAD.

E Details of MR_{KDE}

E.1 Risk with the joint probability density function

Since the conditional and joint densities yield the same maximizer with respect to the output variable, the modal regression function f_M can be defined from the joint density $p_{y\mathbf{x}}(y, \mathbf{x})$ as

$$f_M(\mathbf{x}) = \operatorname{argmax}_t p_{y|\mathbf{x}}(t|\mathbf{x}) = \operatorname{argmax}_t p_{y\mathbf{x}}(t, \mathbf{x}). \quad (42)$$

Thus, the following risk alternative to the modal regression risk can be used for modal regression:

$$\mathcal{R}_J(f) := \int p_{y\mathbf{x}}(f(\mathbf{x}), \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

The following inequality, which follows from (42), ensures that the maximizer of $\mathcal{R}_J(f)$ is f_M :

$$\mathcal{R}_J(f) \leq \int p_{y\mathbf{x}}(f_M(\mathbf{x}), \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

¹The same strategy has been used in Debruyne et al. [2008] that parameter estimation and model selection are performed by the Huber and LAD loss, respectively.

With a parametrized model $f_{\boldsymbol{\theta}}(\mathbf{x})$ as in the kernel model, the empirical version of \mathcal{R}_J can be obtained as

$$\widehat{\mathcal{R}}_J(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n p_{y\mathbf{x}}(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{x}_i).$$

In practice, we need to estimate the joint density $p_{y\mathbf{x}}(y, \mathbf{x})$ to approximate $\widehat{\mathcal{R}}_J(\boldsymbol{\theta})$. Below, we employ kernel density estimation (KDE) for the joint density $p_{y\mathbf{x}}(y, \mathbf{x})$ as done in Yao et al. [2012] and derive an update rule similar as DMR-K.

E.2 Update rule based on a fixed-point method

Let us define KDE with the Gaussian kernel for joint density estimation by

$$\widehat{p}_{\text{KDE}}(y, \mathbf{x}) = \frac{1}{nZ} \sum_{l=1}^n \exp\left(-\frac{(y - y_l)^2}{2h_y^2}\right) \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_l\|^2}{2h_x^2}\right),$$

where $Z = (2\pi)^{(d_x+1)/2} h_y h_x^{d_x}$, and h_y and h_x are positive width parameters. Then, $\widehat{p}_{\text{KDE}}(y, \mathbf{x})$ enables us to approximate $\widehat{\mathcal{R}}_J(\boldsymbol{\theta})$ as

$$\widetilde{\mathcal{R}}_{\text{KDE}}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \widehat{p}_{\text{KDE}}(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{x}_i).$$

Computing the gradient of $\widetilde{\mathcal{R}}_{\text{KDE}}(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ yields

$$\frac{\partial}{\partial \boldsymbol{\theta}} \widetilde{\mathcal{R}}_{\text{KDE}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\mathbf{x}_i) \frac{\partial}{\partial y} \widehat{p}_{\text{KDE}}(y, \mathbf{x}) = \frac{1}{n^2 h_y^2 Z} \{ \mathbf{h}_{\text{KDE}}(\boldsymbol{\theta}) - \mathbf{H}_{\text{KDE}}(\boldsymbol{\theta}) \boldsymbol{\theta} \}, \quad (43)$$

where $f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{k}_m(\mathbf{x})$,

$$\begin{aligned} \mathbf{H}_{\text{KDE}}(\boldsymbol{\theta}) &= \sum_{i=1}^n \sum_{l=1}^n \exp\left(-\frac{(\boldsymbol{\theta}^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2h_y^2}\right) \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_l\|^2}{2h_x^2}\right) \mathbf{k}_m(\mathbf{x}_i) \mathbf{k}_m(\mathbf{x}_l)^\top, \\ \mathbf{h}_{\text{KDE}}(\boldsymbol{\theta}) &= \sum_{i=1}^n \sum_{l=1}^n y_l \exp\left(-\frac{(\boldsymbol{\theta}^\top \mathbf{k}_m(\mathbf{x}_i) - y_l)^2}{2h_y^2}\right) \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_l\|^2}{2h_x^2}\right) \mathbf{k}_m(\mathbf{x}_i). \end{aligned}$$

Setting the right-hand side on (43) to equal to zero leads to the following update rule:

$$\boldsymbol{\theta} \leftarrow \mathbf{H}_{\text{KDE}}^{-1}(\boldsymbol{\theta}) \mathbf{h}_{\text{KDE}}(\boldsymbol{\theta}). \quad (44)$$

Eq.(44) is iteratively used to update $\boldsymbol{\theta}$ as in Algorithm 1.

F Other results on artificial data

In addition to (M1) and (M2), we performed the same experiments as Section 5.1 under the following true modal regression function:

$$(M3) \quad f^*(\mathbf{x}) = \frac{1}{d_x} \sum_{j=1}^{d_x} x^{(j)}.$$

Fig.1 plots estimates of f^* by all methods in $d_x = 1$ over all types of noises. Regarding the Gaussian noise, all methods perform well. LAD and Huber give better estimates than KRR and VHGP for the outlier noise, but does not work to the skewed noise. This would be because LAD and Huber assume symmetric noises. Overall, MR_{KDE} and DMR-K perform fairly well to all types of noises. Table 3 is the results for (M3) over various data dimensions and noises, and shows the superior performance of DMR-K on a wide-range of noises and data dimensions as already demonstrated in the main text.

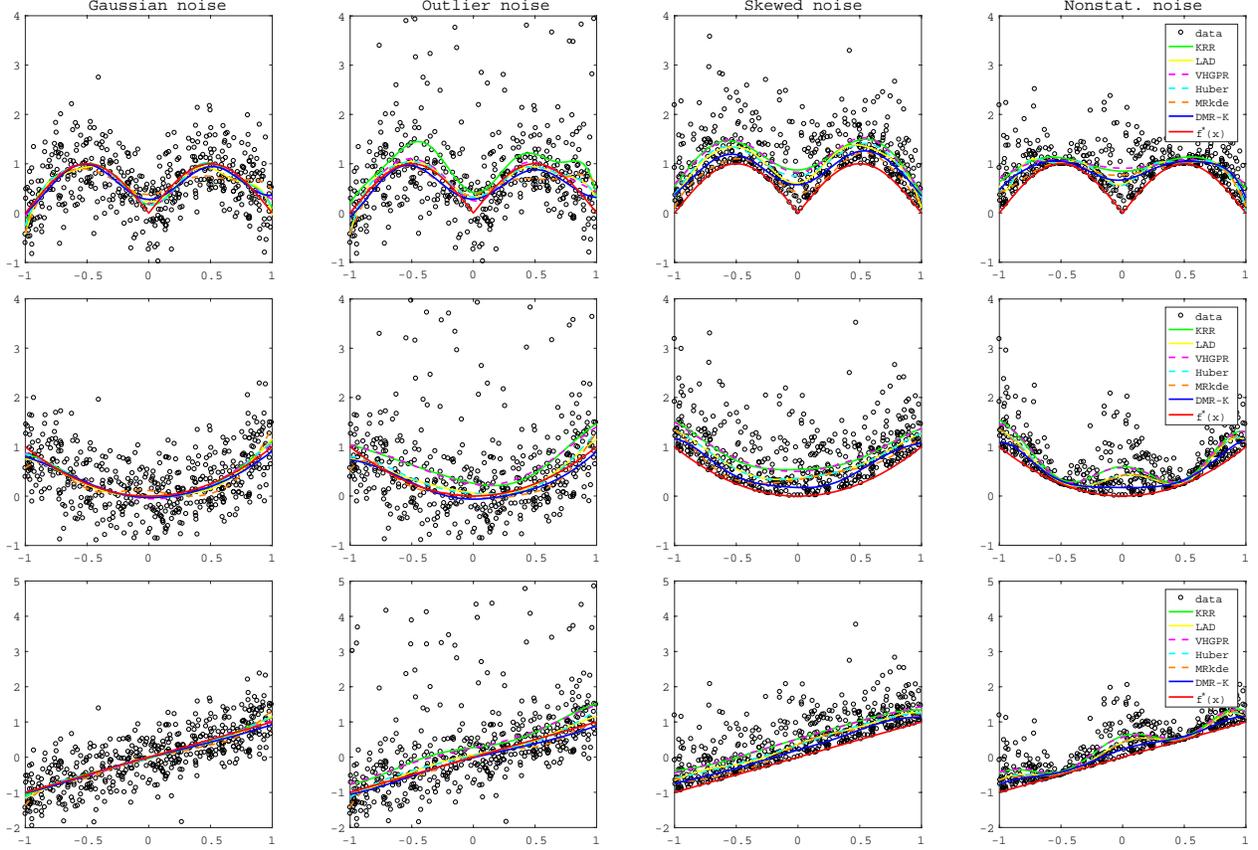


Figure 1: Estimates of $f^*(x)$. The top, middle and the bottom rows are the plots when f^* is (M1), (M2), and (M3), respectively.

G Validity of the performance score (27)

Feng et al. [2017] indicated that the meaning of the maximizer of the surrogate risk $\widehat{\mathcal{R}}^\sigma$ changes depending on the width parameter σ as follows [Feng et al., 2017, Table 2]:

- The maximizer is a conditional *mode* estimator as $\sigma \rightarrow 0$ and $n \rightarrow \infty$
- The maximizer is a (robustified) conditional *mean* estimator as $\sigma, n \rightarrow \infty$

In accord with this theory, Table 4 with a large width parameter ($\sigma = 1.0$) shows LS and LAD outperform DMR-NN and DMR-K because these methods estimate the conditional mean and median asymptotically, while DMR-NN often works better than LS and LAD when the width parameter is small (Table 5, $\sigma = 0.01$). Our choice of $\sigma = n_{te}^{-1/5} = (0.2n)^{-1/5}$ in the main text² is in fact a middle of Tables 4 and 5, and approximately $0.09 \leq \sigma \leq 0.27$ among all datasets. Thus, it seems to be a fairly good choice because the standard deviations in Table 5 are often large and the results for too small σ might be unreliable.

References

M. Debruyne, M. Hubert, and J. A. Suykens. Model selection in kernel based regression using the influence function. *Journal of Machine Learning Research*, 9:2377–2400, 2008.

Y. Feng, J. Fan, and J. A. Suykens. A statistical learning approach to modal regression. *arXiv:1702.05960*, 2017.

²Let us remind that we used 20% of data samples for test in experiments on benchmark datasets (i.e., $n_{te} = 0.2n$).

Table 3: Averages of estimation errors over 30 runs for (M3). The numbers in parentheses indicate standard deviations. The best and comparable methods judged by the t-test at the significance level 1% are described in boldface.

	KRR	LAD	VHGPR	Huber	MR _{KDE}	DMR-K
Gauss noise						
$d_x = 1$	0.04(0.02)	0.06(0.02)	0.03(0.01)	0.05(0.02)	0.10(0.02)	0.05(0.03)
$d_x = 5$	0.07(0.01)	0.09(0.01)	0.05(0.01)	0.07(0.01)	0.19(0.02)	0.06(0.02)
$d_x = 10$	0.08(0.01)	0.11(0.02)	0.10(0.02)	0.08(0.02)	0.29(0.07)	0.08(0.04)
Outlier noise						
$d_x = 1$	0.31(0.03)	0.09(0.02)	0.19(0.11)	0.09(0.02)	0.10(0.02)	0.06(0.02)
$d_x = 5$	0.30(0.02)	0.11(0.02)	0.25(0.07)	0.11(0.02)	0.20(0.02)	0.08(0.02)
$d_x = 10$	0.31(0.03)	0.13(0.02)	0.16(0.07)	0.13(0.02)	0.29(0.05)	0.09(0.04)
Skewed noise						
$d_x = 1$	0.49(0.02)	0.35(0.02)	0.49(0.02)	0.35(0.02)	0.20(0.05)	0.22(0.01)
$d_x = 5$	0.50(0.02)	0.37(0.03)	0.50(0.02)	0.39(0.03)	0.28(0.02)	0.25(0.02)
$d_x = 10$	0.50(0.02)	0.36(0.03)	0.49(0.03)	0.39(0.03)	0.32(0.03)	0.23(0.02)
Nonstationary noise						
$d_x = 1$	0.31(0.02)	0.22(0.02)	0.31(0.02)	0.22(0.02)	0.17(0.03)	0.15(0.01)
$d_x = 5$	0.32(0.01)	0.20(0.02)	0.31(0.02)	0.20(0.02)	0.16(0.01)	0.15(0.01)
$d_x = 10$	0.32(0.02)	0.20(0.02)	0.31(0.02)	0.21(0.02)	0.19(0.02)	0.14(0.01)

A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.

T. Kanamori, T. Suzuki, and M. Sugiyama. Statistical analysis of kernel-based least-squares density-ratio estimation. *Machine Learning*, 86(3):335–367, 2012.

M. Lázaro-Gredilla and M. K. Titsias. Variational heteroscedastic Gaussian process regression. In *Proceedings of the 28th International Conference on International Conference on Machine Learning (ICML)*, pages 841–848, 2011.

L. Wasserman. *All of nonparametric statistics*. Springer, 2006.

W. Yao, B. G. Lindsay, and R. Li. Local modal regression. *Journal of nonparametric statistics*, 24(3):647–663, 2012.

D. Zhou. Derivative reproducing properties for kernel methods in learning theory. *Journal of Computational and Applied Mathematics*, 220(1-2):456–463, 2008.

Table 4: Averages of the performance score (27) over 20 runs when $\sigma = 1.0$. The numbers in parentheses indicate standard deviations. The best and comparable methods judged by the t-test at the significance level 5% are described in boldface. Note that larger numbers indicate better results.

LS	LAD	DMR-NN	DMR-K
space-ga ($d_x = 6, n = 3107$)			
0.3560(0.0027)	0.3563(0.0031)	0.3574(0.0024)	0.3182(0.0056)
abalone ($d_x = 8, n = 4177$)			
0.3437(0.0029)	0.3449(0.0026)	0.3428(0.0033)	0.3193(0.0051)
cpusmall ($d_x = 12, n = 8192$)			
0.3939(0.0003)	0.3939(0.0003)	0.3941(0.0002)	0.3710(0.0019)
cadata ($d_x = 8, n = 20640$)			
0.3673(0.0006)	0.3683(0.0009)	0.3668(0.0011)	0.3349(0.0052)
energy ($d_x = 24, n = 19735$)			
0.3504(0.0023)	0.3633(0.0014)	0.3594(0.0015)	0.3448(0.0012)
superconductivty ($d_x = 81, n = 21263$)			
0.3844(0.0006)	0.3845(0.0007)	0.3761(0.0027)	0.2654(0.0500)
slice loc. ($d_x = 384, n = 53500$)			
0.3987(0.0001)	0.3988(0.0000)	0.3989(0.0000)	-
sgemm ($d_x = 14, n = 241600$)			
0.3982(0.0001)	0.3980(0.0001)	0.3980(0.0001)	0.3726(0.0006)
yearpred. ($d_x = 90, n = 515345$)			
0.3305(0.0006)	0.3361(0.0004)	0.3160(0.0181)	0.1904(0.1441)

Table 5: Averages of the performance score (27) over 20 runs when $\sigma = 0.01$.

LS	LAD	DMR-NN	DMR-K
space-ga ($d_x = 6, n = 3107$)			
0.8118(0.1974)	0.9480(0.1743)	0.9678(0.2138)	0.6135(0.1857)
abalone ($d_x = 8, n = 4177$)			
0.8727(0.1594)	0.9133(0.2500)	0.9544(0.1899)	0.6874(0.1889)
cpusmall ($d_x = 12, n = 8192$)			
3.4789(0.2821)	3.8954(0.3286)	3.5118(0.3599)	2.7954(0.2835)
cadata ($d_x = 8, n = 20640$)			
1.3059(0.0766)	1.5643(0.0862)	1.6245(0.1159)	0.7764(0.0751)
energy ($d_x = 24, n = 19735$)			
1.1987(0.1110)	2.6351(0.1918)	2.8134(0.1349)	2.2965(0.2029)
superconductivty ($d_x = 81, n = 21263$)			
3.3490(0.2462)	5.0061(0.3092)	5.2782(0.4759)	1.3480(0.3183)
slice loc. ($d_x = 384, n = 53500$)			
14.5556(1.0671)	20.7024(0.5988)	24.5490(1.2943)	-
sgemm ($d_x = 14, n = 241600$)			
10.8187(0.7424)	14.0139(0.9916)	12.7494(0.7955)	1.9253(0.1047)
yearpred. ($d_x = 90, n = 515345$)			
0.7629(0.0242)	0.9252(0.0196)	0.9348(0.0926)	0.3590(0.3075)