
Supplementary material for “Robust contrastive learning and nonlinear ICA in the presence of outliers”

Hiroaki Sasaki¹, Takashi Takenouchi^{1,2}, Ricardo Monti³, Aapo Hyvärinen^{4,5}

¹Future University Hakodate, Hokkaido, Japan ²RIKEN AIP, Tokyo, Japan ³University College London, UK

⁴Université Paris-Saclay, Inria, CEA, France ⁵University of Helsinki, Finland

Acknowledgement

The authors would like to thank Dr. Hiroshi Morioka for sharing his PCL codes with us. H.S. was supported by JSPS KAKENHI Grant Number 18K18107. T.T. was supported by JSPS KAKENHI Grant Number 20K03753. A.H. was supported by a Fellowship from CIFAR, and by the DATAIA convergence institute as part of the “Programme d’Investissement d’Avenir”, (ANR-17-CONV-0003) operated by Inria.

A Proof of Theorem 1

Proof. Let us express the inverse of \mathbf{f} in the data generative model (1) by \mathbf{g} such that $\mathbf{s} = \mathbf{g}(\mathbf{x})$. The change of variables provides

$$\begin{aligned} \log p(\mathbf{x}|\mathbf{u}) &= \log \{ (1 - \epsilon(\mathbf{u}))p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}) + \epsilon(\mathbf{u})\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}) \} + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})| \\ &= \log \{ p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}) + \epsilon(\mathbf{u})(\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}) - p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})) \} + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})| \\ &= \log \left[p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}) \left\{ 1 + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u})}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})} - 1 \right) \right\} \right] + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})| \\ &= \log p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}) + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u})}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})} - 1 \right) + O(\epsilon(\mathbf{u})^2) + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})|, \end{aligned}$$

where we applied $\log(1 + \epsilon(\mathbf{u})z) = \epsilon(\mathbf{u})z + O(\epsilon(\mathbf{u})^2)$ with a sufficiently small $\epsilon(\mathbf{u})$ on the last line. Then, the conditionally exponential family assumption (A2) gives

$$\log p(\mathbf{x}|\mathbf{u}) = \sum_{j=1}^{d_x} \lambda_j(\mathbf{u})q_j^*(g_j(\mathbf{x})) + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u})}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})} - 1 \right) + O(\epsilon(\mathbf{u})^2) + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})| - \log Z(\boldsymbol{\lambda}(\mathbf{u})).$$

Contrasting two log-conditional densities of \mathbf{x} given \mathbf{u} and a fixed point \mathbf{u}_0 yields

$$\begin{aligned} \log p(\mathbf{x}|\mathbf{u}) - \log p(\mathbf{x}|\mathbf{u}_0) &= (\boldsymbol{\lambda}(\mathbf{u}) - \boldsymbol{\lambda}(\mathbf{u}_0))^\top \mathbf{q}^*(\mathbf{g}(\mathbf{x})) + \left\{ \epsilon(\mathbf{u}) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u})}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})} - \epsilon(\mathbf{u}_0) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)} \right\} \\ &\quad + O(\epsilon(\mathbf{u})^2) + O(\epsilon(\mathbf{u}_0)^2) - (\log Z(\boldsymbol{\lambda}(\mathbf{u})) - \log Z(\boldsymbol{\lambda}(\mathbf{u}_0))), \end{aligned} \tag{21}$$

where $\boldsymbol{\lambda}(\mathbf{u}) := (\lambda_1(\mathbf{u}), \dots, \lambda_{d_x}(\mathbf{u}))^\top$, $\mathbf{q}^*(\mathbf{g}(\mathbf{x})) := (q_1^*(g_1(\mathbf{x})), \dots, q_{d_x}^*(g_{d_x}(\mathbf{x})))^\top$, and note that the Jacobian $|\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})|$ is cancelled out.

On the other hand, by the universal approximation assumption (A4), we obtain

$$\log p(\mathbf{x}|\mathbf{u}) - \log p(\mathbf{x}|\mathbf{u}_0) = (\mathbf{w}(\mathbf{u}) - \mathbf{w}(\mathbf{u}_0))^\top \mathbf{h}(\mathbf{x}) - \log e(\mathbf{u}) + \log e(\mathbf{u}_0). \tag{22}$$

Then, equating (21) with (22) provides

$$\bar{\lambda}(\mathbf{u})^\top \mathbf{q}^*(\mathbf{g}(\mathbf{x})) + \epsilon(\mathbf{u}) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u})}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u})} - \epsilon(\mathbf{u}_0) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)} + O(\epsilon(\mathbf{u})^2) + O(\epsilon(\mathbf{u}_0)^2) = \bar{\mathbf{w}}(\mathbf{u})^\top \mathbf{h}(\mathbf{x}) + \bar{\beta}(\mathbf{u}), \quad (23)$$

where $\bar{\mathbf{w}}(\mathbf{u}) := \mathbf{w}(\mathbf{u}) - \mathbf{w}(\mathbf{u}_0)$, $\bar{\lambda}(\mathbf{u}) := \lambda(\mathbf{u}) - \lambda(\mathbf{u}_0)$, and $\bar{\beta}(\mathbf{u}) := \log e(\mathbf{u}) - \log e(\mathbf{u}_0) + \log Z(\lambda(\mathbf{u})) - \log Z(\lambda(\mathbf{u}_0))$.

Next, we multiply $\bar{\lambda}(\mathbf{u})$ to the both sides of (23) and evaluate it at m points, $\mathbf{u}_1, \dots, \mathbf{u}_m$. Finally, taking the summation for $\mathbf{u}_1, \dots, \mathbf{u}_m$ yields

$$\underbrace{\left(\sum_{i=1}^m \bar{\lambda}(\mathbf{u}_i) \bar{\lambda}(\mathbf{u}_i)^\top \right)}_{\bar{\Lambda}} \mathbf{q}^*(s) + \sum_{i=1}^m \left(\epsilon(\mathbf{u}_i) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}_i)}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}_i)} - \epsilon(\mathbf{u}_0) \frac{\delta(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)}{p^*(\mathbf{g}(\mathbf{x})|\mathbf{u}_0)} \right) \bar{\lambda}(\mathbf{u}_i) + \sum_{i=1}^m \{O(\epsilon(\mathbf{u}_i)^2) + \epsilon(\mathbf{u}_0)^2\} \bar{\lambda}(\mathbf{u}_i) = \left(\sum_{i=1}^m \bar{\mathbf{w}}(\mathbf{u}_i) \bar{\lambda}(\mathbf{u}_i)^\top \right) \mathbf{h}(\mathbf{x}) + \sum_{i=1}^m \bar{\beta}(\mathbf{u}_i) \bar{\lambda}(\mathbf{u}_i).$$

Applying the inverse of $\bar{\Lambda}$ to both sides completes the proof. \square

B Influence of outliers in general non-exponential case

As in Hyvärinen et al. [2019], we assume that the general conditional independence (4) holds, and that both the mixing function \mathbf{f} in the data generative model (1) and a nonlinear feature $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_{d_x}(\mathbf{x}))^\top$ are invertible. Then, by the change of variables from the generative model (1),

$$\mathbf{v}(\mathbf{y}) := \mathbf{s} = \mathbf{f}^{-1} \circ \mathbf{h}^{-1}(\mathbf{y}), \text{ and } \mathbf{y} := \mathbf{h}(\mathbf{x})$$

where \circ denotes composition. For the case of no outliers, Hyvärinen et al. [2019] proved that the nonlinear feature $\mathbf{h}(\mathbf{x})$ asymptotically recover the source \mathbf{s} up to elementwise invertible transformations by showing that the following equations hold: For all $i, j, k = 1, \dots, d_x$ with $j \neq k$,

$$\frac{\partial}{\partial y_j} v_i(\mathbf{y}) = 0, \quad \frac{\partial}{\partial y_k} v_i(\mathbf{y}) = 0, \quad \text{and} \quad \frac{\partial^2}{\partial y_j \partial y_k} v_i(\mathbf{y}) = 0. \quad (24)$$

Eqs.(24) indicate that each $v_i(= s_i)$ is a function of only one distinct element in $\mathbf{y}(= \mathbf{h}(\mathbf{x}))$, and thus \mathbf{s} is identifiable by $\mathbf{h}(\mathbf{x})$ up to elementwise invertible transformations. By denoting the first- and second-order derivatives of $v_i(\mathbf{y})$ in (24) as

$$v_i^j := \frac{\partial}{\partial y_j} v_i(\mathbf{y}), \quad v_i^{j,k} := \frac{\partial^2}{\partial y_j \partial y_k} v_i(\mathbf{y}),$$

Eqs.(24) for all i, j, k can be compactly expressed as the following matrix form:

$$\mathbf{M}(\mathbf{v}) = \mathbf{O}, \quad (25)$$

where \mathbf{O} denotes the null matrix, and $\mathbf{M}(\mathbf{v})$ is a $(d_x^2 - d_x)$ by $2d_x$ matrix and defined by

$$\mathbf{M}(\mathbf{v}) := \begin{pmatrix} v_1^1 v_1^2 & v_2^1 v_2^2 & \cdots & v_{d_x}^1 v_{d_x}^2 & v_1^{12} & v_2^{12} & \cdots & v_{d_x}^{12} \\ v_1^1 v_1^3 & v_2^1 v_2^3 & \cdots & v_{d_x}^1 v_{d_x}^3 & v_1^{13} & v_2^{13} & \cdots & v_{d_x}^{13} \\ \vdots & \vdots & \cdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ v_1^{d_x-1} v_1^{d_x} & v_2^{d_x-1} v_2^{d_x} & \cdots & v_{d_x}^{d_x-1} v_{d_x}^{d_x} & v_1^{d_x-1, d_x} & v_2^{d_x-1, d_x} & \cdots & v_{d_x}^{d_x-1, d_x} \end{pmatrix}.$$

However, in the presence of outliers, (25) does not hold in general. In order to investigate the influence from outliers, we define the following notations:

$$\begin{aligned}\widetilde{\mathbf{M}}(\mathbf{v}) &:= - \begin{pmatrix} v_1^1 v_2^2 & v_1^1 v_3^2 & \cdots & v_{d_x-1}^1 v_{d_x}^2 \\ v_1^1 v_2^3 & v_1^1 v_3^3 & \cdots & v_{d_x-1}^1 v_{d_x}^3 \\ \vdots & \vdots & \cdots & \vdots \\ v_1^{d_x-1} v_2^{d_x} & v_1^{d_x-1} v_3^{d_x} & \cdots & v_{d_x-1}^{d_x-1} v_{d_x}^{d_x} \end{pmatrix} \\ \mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) &:= \left(\frac{\partial^2}{\partial v_1^2} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0), \dots, \frac{\partial^2}{\partial v_{d_x}^2} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0), \frac{\partial}{\partial v_1} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0), \dots, \frac{\partial}{\partial v_{d_x}} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0) \right)^\top \\ \widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) &:= \left(\frac{\partial^2}{\partial v_1 \partial v_2} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0), \frac{\partial^2}{\partial v_1 \partial v_3} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0), \dots, \frac{\partial^2}{\partial v_{d_x-1} \partial v_{d_x}} \varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0) \right)^\top,\end{aligned}$$

where $\varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0) := \log p(\mathbf{v}|\mathbf{u}) - \log p(\mathbf{v}|\mathbf{u}_0)$, \mathbf{u}_0 denotes a fixed point of \mathbf{u} , $\widetilde{\mathbf{M}}(\mathbf{v})$ by a $(d_x^2 - d_x)$ by $(d_x^2 - d_x)$ matrix, while $\mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ and $\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ is the $2d_x$ - and $(d_x^2 - d_x)$ -dimensional vectors, respectively. The following theorem is useful to understand how (25) is affected by outliers:

Theorem 3. Assume that

(B1) Data \mathbf{x} is generated from (1) where \mathbf{f} is invertible.

(B2) Conditional independence (4) holds.

(B3) For all \mathbf{v} and \mathbf{u} , $\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$, $\frac{\delta^i(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$ and $\frac{\delta^{i,j}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$ are finite where $\delta^i(\mathbf{v}|\mathbf{u}) := \frac{\partial}{\partial v_i} \delta(\mathbf{v}|\mathbf{u})$ and $\delta^{i,j}(\mathbf{v}|\mathbf{u}) := \frac{\partial^2}{\partial v_i \partial v_j} \delta(\mathbf{v}|\mathbf{u})$.

(B4) The conditional density of \mathbf{x} given \mathbf{u} is universally approximated with an invertible feature extractor $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_{d_x}(\mathbf{x}))$ as

$$\log \frac{p(\mathbf{x}|\mathbf{u})}{c(\mathbf{x})e(\mathbf{u})} = \sum_{i=1}^{d_x} \psi(h_i(\mathbf{x}), \mathbf{u}), \quad (26)$$

where ψ , c and e are nonzero scalar functions.

Then, under the contaminated density model (5), the following holds at a fixed point \mathbf{u}_0 :

$$\mathbf{M}(\mathbf{v})\mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) = \widetilde{\mathbf{M}}(\mathbf{v})\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0). \quad (27)$$

The proof is given in Section B.1. First of all, in the case of no outliers, Theorem 3 essentially recovers Theorem 1 in Hyvärinen et al. [2019]: When $\epsilon(\mathbf{u}) = 0$ for all \mathbf{u} , $\varphi(\mathbf{v}|\mathbf{u}, \mathbf{u}_0) = \sum_{i=1}^{d_x} \{q^*(v_i|\mathbf{u}_0) - q^*(v_i|\mathbf{u})\}$ and thus $\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) = \mathbf{0}$, which yields from (27)

$$\mathbf{M}(\mathbf{v})\mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) = \mathbf{O}.$$

Then, applying an additional assumption called *Assumption of Variability* [Hyvärinen et al., 2019]¹ recovers (25).

However, when $\epsilon(\mathbf{u}) \neq 0$, $\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ is a nonzero vector and thus the right-hand sides on (27) are nonzeros in general. Let us remind that the elements in $\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ are given by

$$\frac{\partial^2}{\partial v_i \partial v_j} \{\log p(\mathbf{v}|\mathbf{u}) - \log p(\mathbf{v}|\mathbf{u}_0)\}, \quad (28)$$

where $i, j = 1, \dots, d_x$ but $i \neq j$. Obviously, it is not easy to understand when these elements (28) are strongly deviated from zeros, yet the following proposition proved in Section B.2 gives some insight:

¹The assumption mean that there exist $2d_x + 1$ points, $\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{2d_x}$ such that $\mathbf{w}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_0), \dots, \mathbf{w}(\mathbf{v}, \mathbf{u}_{2d_x}, \mathbf{u}_0)$ are linear independent, implying that the conditional density $p^*(\mathbf{v}|\mathbf{u})$ is sufficiently complex and diverse.

Proposition 1. Assume that the conditional independence (4) holds and sources are generated from the contaminated density model (5). Then, with sufficiently small $\epsilon(\mathbf{u})$,

$$\begin{aligned} & \frac{\partial^2}{\partial v_i \partial v_j} \log p(\mathbf{v}|\mathbf{u}) \\ &= \epsilon(\mathbf{u}) \left\{ \frac{\delta^{i,j}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - q^{*i}(v_i|\mathbf{u}) \frac{\delta^j(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - q^{*j}(v_j|\mathbf{u}) \frac{\delta^i(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} + q^{*i}(v_i|\mathbf{u}) q^{*j}(v_j|\mathbf{u}) \frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \right\} + O(\epsilon(\mathbf{u})^2), \end{aligned} \quad (29)$$

where $i \neq j$ and $q^{*i}(v_i|\mathbf{u}) := \frac{\partial}{\partial v_i} q^*(v_i|\mathbf{u})$.

Eq.(29) allows us to more easily understand the implication of (28): $\tilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ would be strongly deviated from the zero vector when at least one of $\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$, $\frac{\delta^i(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$, $\frac{\delta^j(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$ and $\frac{\delta^{i,j}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})}$ is large. The four ratio factors could be very large when smooth $\delta(\mathbf{v}|\mathbf{u})$ lies on the tails of $p^*(\mathbf{v}|\mathbf{u})$.

B.1 Proof of Theorem 3

Proof. We first obtain the expression of $\log p(\mathbf{x}|\mathbf{u})$ by the change of variables the data generative model (1) as

$$\log p(\mathbf{x}|\mathbf{u}) = \log p(\mathbf{g}(\mathbf{x})|\mathbf{u}) - \log Z(\mathbf{u}) + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})|,$$

where $\mathbf{g} := \mathbf{f}^{-1}$, and $\mathbf{J}_{\mathbf{g}}(\mathbf{x})$ and \det denote the Jacobian and determinant, respectively. Then, the universal approximation assumption (B4) gives

$$\sum_{i=1}^{d_x} \psi(h_i(\mathbf{x}), \mathbf{u}) = \log p(\mathbf{g}(\mathbf{x})|\mathbf{u}) - \log Z(\mathbf{u}) + \log |\det \mathbf{J}_{\mathbf{g}}(\mathbf{x})|.$$

To remove the Jacobian term and log-partition function, we compute the differences of the above equations at \mathbf{u} and a fixed point \mathbf{u}_0 as

$$\sum_{i=1}^{d_x} \bar{\psi}(h_i(\mathbf{x}), \mathbf{u}, \mathbf{u}_0) = \log p(\mathbf{g}(\mathbf{x})|\mathbf{u}) - \log p(\mathbf{g}(\mathbf{x})|\mathbf{u}_0).$$

where

$$\bar{\psi}(h_i(\mathbf{x}), \mathbf{u}, \mathbf{u}_0) := \psi(h_i(\mathbf{x}), \mathbf{u}) - \psi(h_i(\mathbf{x}), \mathbf{u}_0).$$

By the further change of variables $\mathbf{y} = \mathbf{h}(\mathbf{x})$ and $\mathbf{v}(\mathbf{y}) = \mathbf{g}(\mathbf{h}^{-1}(\mathbf{y}))$,

$$\sum_{i=1}^{d_x} \bar{\psi}(y_i, \mathbf{u}, \mathbf{u}_0) = \varphi(\mathbf{v}(\mathbf{y}), \mathbf{u}, \mathbf{u}_0), \quad (30)$$

where

$$\varphi(\mathbf{v}(\mathbf{y}), \mathbf{u}, \mathbf{u}_0) := \log p(\mathbf{v}(\mathbf{y})|\mathbf{u}) - \log p(\mathbf{v}(\mathbf{y})|\mathbf{u}_0).$$

Next, let us use or remind the following notations:

$$\begin{aligned} \varphi^l(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) &:= \frac{\partial}{\partial v_l} \varphi(\mathbf{v}, \mathbf{u}, \mathbf{u}_0), & \varphi^{l,m}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) &:= \frac{\partial^2}{\partial v_l \partial v_m} \varphi(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) \\ v_l^k &:= \frac{\partial}{\partial y_k} v_l(\mathbf{y}), & v_l^{j,k} &:= \frac{\partial^2}{\partial y_j \partial y_k} v_l(\mathbf{y}). \end{aligned}$$

Taking the second-order partial derivative of the left-hand side on (30) with respect to y_j and y_k for $j \neq k$ yields

$$\frac{\partial^2}{\partial y_j \partial y_k} \sum_{i=1}^{d_x} \bar{\psi}(y_i, \mathbf{u}, \mathbf{u}_0) = 0. \quad (31)$$

On the other hand, the second-order partial derivative of the right-hand side on (30) can be expressed as

$$\frac{\partial^2}{\partial y_j \partial y_k} \varphi(\mathbf{v}(\mathbf{y}), \mathbf{u}, \mathbf{u}_0) = \sum_{l=1}^{d_x} \left[\varphi^{l,l}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^j v_l^k + \varphi^l(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^{jk} \right] + \sum_{l=1}^{d_x} \sum_{\substack{m=1 \\ m \neq l}}^{d_x} \varphi^{l,m}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^j v_m^k. \quad (32)$$

Equating (31) with (32) under (30) gives the following equation:

$$\sum_{l=1}^{d_x} \left[\varphi^{l,l}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^j v_l^k + \varphi^l(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^{jk} \right] = - \sum_{l=1}^{d_x} \sum_{\substack{m=1 \\ m \neq l}}^{d_x} \varphi^{l,m}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) v_l^j v_m^k. \quad (33)$$

Regarding $j = 1, \dots, d_x$ and $k = 1, \dots, d_x$, we collect all of $v_l^j v_l^k$, v_l^{jk} and $-v_l^j v_m^k$ for $j \neq k$ and $l \neq m$ as $(d_x^2 - d_x)$ -dimensional vectors, $\mathbf{a}_l(\mathbf{v})$, $\mathbf{b}_l(\mathbf{v})$, and $\mathbf{c}_{l,m}(\mathbf{v})$, respectively. Then, (33) for all j and k (but $j \neq k$) can be expressed as

$$\sum_{l=1}^{d_x} \left[\varphi^{l,l}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) \mathbf{a}_l(\mathbf{v}) + \varphi^l(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) \mathbf{b}_l(\mathbf{v}) \right] = \sum_{l=1}^{d_x} \sum_{\substack{m=1 \\ m \neq l}}^{d_x} \varphi^{l,m}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) \mathbf{c}_{l,m}(\mathbf{v}).$$

Furthermore, the above equations for all l and m (but $l \neq m$) can be summarized as the following system of linear equations:

$$\mathbf{M}(\mathbf{v}) \mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0) = \widetilde{\mathbf{M}}(\mathbf{v}) \widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0),$$

where $\mathbf{M}(\mathbf{v}) := (\mathbf{a}_1(\mathbf{v}), \dots, \mathbf{a}_n(\mathbf{v}), \mathbf{b}_1(\mathbf{v}), \dots, \mathbf{b}_n(\mathbf{v}))$, $\widetilde{\mathbf{M}}(\mathbf{v}) := (\mathbf{c}_{1,2}(\mathbf{v}), \mathbf{c}_{1,3}(\mathbf{v}), \dots, \mathbf{c}_{d_x, d_x-1}(\mathbf{v}))$, and $\mathbf{w}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ is the $2d_x$ -dimensional vector of all $\varphi^{l,l}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ and $\varphi^l(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ for $l = 1, \dots, d_x$, and $\widetilde{\mathbf{w}}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ is the $(d_x^2 - d_x)$ -dimensional vector of all $\varphi^{l,m}(\mathbf{v}, \mathbf{u}, \mathbf{u}_0)$ for $l, m = 1, \dots, d_x$ but $l \neq m$. Thus, the proof is completed. \square

B.2 Proof of Proposition 1

Proof. We first define the following notations:

$$\begin{aligned} p^l(\mathbf{v}|\mathbf{u}) &:= \frac{\partial}{\partial v_l} p(\mathbf{v}|\mathbf{u}), & p^{l,m}(\mathbf{v}|\mathbf{u}) &:= \frac{\partial^2}{\partial v_l \partial v_m} p(\mathbf{v}|\mathbf{u}) \\ p^{*l}(\mathbf{v}|\mathbf{u}) &:= \frac{\partial}{\partial v_l} p^*(\mathbf{v}|\mathbf{u}), & p^{*l,m}(\mathbf{v}|\mathbf{u}) &:= \frac{\partial^2}{\partial v_l \partial v_m} p^*(\mathbf{v}|\mathbf{u}). \end{aligned}$$

Then, under the contaminated density model (5), we compute

$$\begin{aligned} \frac{\partial^2}{\partial v_l \partial v_m} \log p(\mathbf{v}|\mathbf{u}) &= \frac{\partial^2}{\partial v_l \partial v_m} \log p^*(\mathbf{v}|\mathbf{u}) + \frac{\partial^2}{\partial v_l \partial v_m} \log \left\{ 1 + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - 1 \right) \right\} \\ &= \frac{\partial^2}{\partial v_l \partial v_m} \log \left\{ 1 + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - 1 \right) \right\}, \end{aligned} \quad (34)$$

where we used the following relation under the conditional independence assumption (4):

$$\frac{\partial^2}{\partial v_l \partial v_m} \log p^*(\mathbf{v}|\mathbf{u}) = \frac{\partial^2}{\partial v_l \partial v_m} \sum_{i=1}^{d_x} q^*(v_i|\mathbf{u}) = 0.$$

Next, we apply the Taylor expansion with sufficiently small $\epsilon(\mathbf{u})$ to the right-hand side on (34):

$$\begin{aligned} &\frac{\partial^2}{\partial v_l \partial v_m} \log \left\{ 1 + \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - 1 \right) \right\} \\ &= \frac{\partial^2}{\partial v_l \partial v_m} \left\{ \epsilon(\mathbf{u}) \left(\frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - 1 \right) + O(\epsilon(\mathbf{u})^2) \right\} \\ &= \epsilon(\mathbf{u}) \left\{ \frac{\delta^{l,m}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - \frac{\delta^l(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \frac{p^{*m}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - \frac{\delta^m(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \frac{p^{*l}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} + \frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \frac{p^{*l}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \frac{p^{*n}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \right\} + O(\epsilon(\mathbf{u})^2). \end{aligned} \quad (35)$$

Substituting (35) into (34) yields

$$\begin{aligned} & \frac{\partial^2}{\partial v_l \partial v_m} \log p(\mathbf{v}|\mathbf{u}) \\ &= \epsilon(\mathbf{u}) \left\{ \frac{\delta^{l,m}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - q^{*l}(v_l|\mathbf{u}) \frac{\delta^m(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} - q^{*m}(v_m|\mathbf{u}) \frac{\delta^l(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} + q^{*l}(v_l|\mathbf{u}) q^{*m}(v_m|\mathbf{u}) \frac{\delta(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} \right\} + O(\epsilon(\mathbf{u})^2), \end{aligned}$$

where we applied

$$\frac{p^{*l}(\mathbf{v}|\mathbf{u})}{p^*(\mathbf{v}|\mathbf{u})} = \frac{\partial}{\partial v_l} \log p^*(\mathbf{v}|\mathbf{u}) = q^{*l}(v_l|\mathbf{u}),$$

under the conditional independence assumption (4). Thus, the proof is completed. \square

C Proof of Theorem 2

C.1 Derivation of (13)

With $p(y, \mathbf{x}, \mathbf{u}) = p(\mathbf{x}, \mathbf{u}|y)p(y)$, we have

$$\begin{aligned} d_\gamma(p(y|\mathbf{x}, \mathbf{u}), r(\mathbf{x}, \mathbf{u}); p(\mathbf{x}, \mathbf{u})) &:= -\frac{1}{\gamma} \log \iint \sum_{y=0}^1 \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{y(\gamma+1)}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(y, \mathbf{x}, \mathbf{u}) d\mathbf{x}d\mathbf{u} \\ &= -\frac{1}{\gamma} \log \left[\iint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x}, \mathbf{u}|y=1)p(y=1) d\mathbf{x}d\mathbf{u} \right. \\ &\quad \left. + \iint \left\{ \frac{1}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x}, \mathbf{u}|y=0)p(y=0) d\mathbf{x}d\mathbf{u} \right] \\ &= -\frac{1}{\gamma} \log \left[\underbrace{\frac{1}{2} \iint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x}, \mathbf{u}) d\mathbf{x}d\mathbf{u}}_{(A)} \right. \\ &\quad \left. + \frac{1}{2} \iint \left\{ \frac{1}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x})p(\mathbf{u}) d\mathbf{x}d\mathbf{u} \right], \end{aligned} \quad (36)$$

where $p(y=0) = p(y=1) = \frac{1}{2}$, $p(\mathbf{x}, \mathbf{u}|y=1) = p(\mathbf{x}, \mathbf{u})$, and $p(\mathbf{x}, \mathbf{u}|y=0) = p(\mathbf{x})p(\mathbf{u})$. Under the outlier model, the joint density can be expressed as

$$p(\mathbf{x}, \mathbf{u}) = p(\mathbf{x}|\mathbf{u})p(\mathbf{u}) = (1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}) + \epsilon(\mathbf{u})\delta(\mathbf{x}, \mathbf{u}).$$

Then, the term (A) in (36) can be written as

$$\begin{aligned} (A) &= \iint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x}, \mathbf{u}) d\mathbf{x}d\mathbf{u} \\ &= \iint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} (1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}) d\mathbf{x}d\mathbf{u} + \iint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \epsilon(\mathbf{u})\delta(\mathbf{x}, \mathbf{u}) d\mathbf{x}d\mathbf{u} \end{aligned} \quad (37)$$

Finally, substituting (37) into (36) yields

$$\begin{aligned}
& d_\gamma(p(y|\mathbf{x}, \mathbf{u}), r(\mathbf{x}, \mathbf{u}); p(\mathbf{x}, \mathbf{u})) \\
&= -\frac{1}{\gamma} \log \left[\frac{1}{2} \iiint \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u})d\mathbf{x}d\mathbf{u} + \frac{1}{2} \iiint \left\{ \frac{1}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} p(\mathbf{x})p(\mathbf{u})d\mathbf{x}d\mathbf{u} \right. \\
&\quad \left. + \frac{1}{2} \iiint \underbrace{\left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1+r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \epsilon(\mathbf{u})\delta(\mathbf{x}, \mathbf{u})d\mathbf{x}d\mathbf{u}}_{=\nu} \right] \\
&= J[r(\mathbf{x}, \mathbf{u}); (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}), p(\mathbf{x})p(\mathbf{u})] + O(\nu),
\end{aligned}$$

where we applied the relation $\log(y + \nu) = \log(y) + O(\nu)$ with sufficiently small ν .

C.2 Proof of the minimizer (14)

C.2.1 Preliminaries

We use the following results in the main proof, which are derived from the Taylor expansion:

$$\left(\frac{1}{(r + \eta\phi)^{\gamma+1} + 1} \right)^{\frac{\gamma}{\gamma+1}} = \left(\frac{1}{r^{\gamma+1} + 1} \right)^{\frac{\gamma}{\gamma+1}} - \eta \frac{\gamma r^\gamma \phi}{(r^{\gamma+1} + 1)^{\frac{2\gamma+1}{\gamma+1}}} + \frac{\eta^2}{2} \frac{\{\gamma(\gamma+1)r^{2\gamma} - \gamma^2 r^{\gamma-1}\} \phi^2}{(r^{\gamma+1} + 1)^{\frac{3\gamma+1}{\gamma+1}}} + O(\eta^3) \quad (38)$$

$$\left(\frac{(r + \eta\phi)^{\gamma+1}}{(r + \eta\phi)^{\gamma+1} + 1} \right)^{\frac{\gamma}{\gamma+1}} = \left(\frac{r^{\gamma+1}}{r^{\gamma+1} + 1} \right)^{\frac{\gamma}{\gamma+1}} + \eta \frac{\gamma r^{\gamma-1} \phi}{(r^{\gamma+1} + 1)^{\frac{2\gamma+1}{\gamma+1}}} - \frac{\eta^2}{2} \frac{\{\gamma(\gamma+2)r^{2\gamma-1} - \gamma(\gamma-1)r^{\gamma-2}\} \phi^2}{(r^{\gamma+1} + 1)^{\frac{3\gamma+2}{\gamma+1}}} + O(\eta^3). \quad (39)$$

C.2.2 Main proof

Proof. Let us define $\tilde{J}[r] := \exp(-\gamma J[r(\mathbf{x}, \mathbf{u}); (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}), p(\mathbf{x})p(\mathbf{u})])$, and then we derive a maximizer of $\tilde{J}[r]$ alternative to a minimizer of $J[r(\mathbf{x}, \mathbf{u}); (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}), p(\mathbf{x})p(\mathbf{u})]$. For $\eta > 0$ and a perturbation ϕ , with (38) and (39), we have

$$\begin{aligned}
\tilde{J}[r + \eta\phi] &= \tilde{J}[r] + \frac{\eta}{2} \iint \left[\frac{\gamma r(\mathbf{x}, \mathbf{u})^{\gamma-1} \phi(\mathbf{x}, \mathbf{u}) \{ (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}) - r(\mathbf{x}, \mathbf{u})p(\mathbf{x})p(\mathbf{u}) \}}{(r(\mathbf{x}, \mathbf{u})^{\gamma+1} + 1)^{\frac{2\gamma+1}{\gamma+1}}} \right] d\mathbf{x}d\mathbf{u} \\
&\quad + \frac{\eta^2}{4} \iint \left[\frac{\{\gamma(\gamma+1)r(\mathbf{x}, \mathbf{u})^{2\gamma} - \gamma^2 r(\mathbf{x}, \mathbf{u})^{\gamma-1}\} p(\mathbf{x})p(\mathbf{u})\phi(\mathbf{x}, \mathbf{u})^2}{(r(\mathbf{x}, \mathbf{u})^{\gamma+1} + 1)^{\frac{3\gamma+2}{\gamma+1}}} \right. \\
&\quad \left. - \frac{\{\gamma(\gamma+2)r(\mathbf{x}, \mathbf{u})^{2\gamma-1} - \gamma(\gamma-1)r(\mathbf{x}, \mathbf{u})^{\gamma-2}\} (1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u})\phi(\mathbf{x}, \mathbf{u})^2}{(r(\mathbf{x}, \mathbf{u})^{\gamma+1} + 1)^{\frac{3\gamma+2}{\gamma+1}}} \right] d\mathbf{x}d\mathbf{u} + O(\eta^3). \quad (40)
\end{aligned}$$

The optimality condition is satisfied when the term of order η on the right-hand side of (40) equals to zero for arbitrary ϕ . Thus, an optimizer $r^*(\mathbf{x}, \mathbf{u})$ satisfies the following equation:

$$(1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}) - r(\mathbf{x}, \mathbf{u})p(\mathbf{x})p(\mathbf{u}) = 0.$$

Then, we obtain $r^*(\mathbf{x}, \mathbf{u})$ as

$$r^*(\mathbf{x}, \mathbf{u}) = \frac{(1-\epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u})}{p(\mathbf{x})p(\mathbf{u})} = \frac{(1-\epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})}{p(\mathbf{x})}.$$

To investigate if r^* is a maximizer of $J[r]$, we compute the term of order η^2 on the right-hand side of (40) at $r = r^*$ as

$$-\iint \left[\frac{\gamma \{ r^*(\mathbf{x}, \mathbf{u})^{2\gamma} + r^*(\mathbf{x}, \mathbf{u})^{\gamma-1} \} p(\mathbf{x})p(\mathbf{u})\phi(\mathbf{x}, \mathbf{u})^2}{(r^*(\mathbf{x}, \mathbf{u})^{\gamma+1} + 1)^{\frac{3\gamma+1}{\gamma+1}}} \right] d\mathbf{x}d\mathbf{u},$$

where we used the relation, $(1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}, \mathbf{u}) = r^*(\mathbf{x}, \mathbf{u})p(\mathbf{x})p(\mathbf{u})$. This shows that the term of order η^2 is negative for any choices of ϕ . Thus, r^* is a maximizer, and the proof is completed. \square

D Proof of Proposition 1

Proof. We first define the support of $\delta(\mathbf{x}|\mathbf{u})$ as

$$\mathcal{X}_{\mathbf{u}}^{\delta} := \{\mathbf{x} \mid \delta(\mathbf{x}|\mathbf{u}) > 0\}.$$

In the neighborhood of $r^*(\mathbf{x}, \mathbf{u}) = \frac{(1-\epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})}{p(\mathbf{x})}$, we can obtain

$$\begin{aligned} \nu &= \iint_{\mathcal{X}, \mathcal{U}} \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1 + r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \epsilon(\mathbf{u}) \delta(\mathbf{x}, \mathbf{u}) d\mathbf{x} d\mathbf{u} \\ &= \int_{\mathcal{U}} \left[\int_{\mathcal{X}_{\mathbf{u}}^{\delta}} \left\{ \frac{r(\mathbf{x}, \mathbf{u})^{\gamma+1}}{1 + r(\mathbf{x}, \mathbf{u})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \epsilon(\mathbf{u}) \delta(\mathbf{x}|\mathbf{u}) d\mathbf{x} \right] p(\mathbf{u}) d\mathbf{u} \\ &\leq \int_{\mathcal{U}} \left[\int_{\mathcal{X}_{\mathbf{u}}^{\delta}} \left\{ \frac{p^*(\mathbf{x}|\mathbf{u})^{\gamma+1}}{\{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})\}^{\gamma+1} + p(\mathbf{x})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \delta(\mathbf{x}|\mathbf{u}) d\mathbf{x} \right] g_{\gamma}(\mathbf{u}) + O\left(\sup_{\mathbf{x}, \mathbf{u}} |r(\mathbf{x}, \mathbf{u}) - r^*(\mathbf{x}, \mathbf{u})|\right), \end{aligned} \quad (41)$$

where $g_{\gamma}(\mathbf{u}) := (1 - \epsilon(\mathbf{u}))^{\gamma} \epsilon(\mathbf{u}) p(\mathbf{u})$. Since the generative model (1) is invertible, by the change of variables from \mathbf{x} to \mathbf{s} , we obtain the following equation:

$$\begin{aligned} &\int_{\mathcal{X}_{\mathbf{u}}^{\delta}} \left\{ \frac{p^*(\mathbf{x}|\mathbf{u})^{\gamma+1}}{\{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})\}^{\gamma+1} + p(\mathbf{x})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \delta(\mathbf{x}|\mathbf{u}) d\mathbf{x} \\ &= \int_{\mathcal{S}_{\mathbf{u}}^{\delta}} \left\{ \frac{p^*(\mathbf{s}|\mathbf{u})^{\gamma+1}}{\{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{s}|\mathbf{u})\}^{\gamma+1} + p(\mathbf{s})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \delta(\mathbf{s}|\mathbf{u}) d\mathbf{s}. \end{aligned} \quad (42)$$

Substituting (42) into (41) yields

$$\nu \leq \int_{\mathcal{U}} \left[\int_{\mathcal{S}_{\mathbf{u}}^{\delta}} \left\{ \frac{p^*(\mathbf{s}|\mathbf{u})^{\gamma+1}}{\{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{s}|\mathbf{u})\}^{\gamma+1} + p(\mathbf{s})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \delta(\mathbf{s}|\mathbf{u}) d\mathbf{s} \right] g_{\gamma}(\mathbf{u}) d\mathbf{u} + O\left(\sup_{\mathbf{x}, \mathbf{u}} |r(\mathbf{x}, \mathbf{u}) - r^*(\mathbf{x}, \mathbf{u})|\right).$$

Under the assumption that $\mathcal{S}_{\mathbf{u}}^{p^*} \cap \mathcal{S}_{\mathbf{u}}^{\delta} = \emptyset$, we can easily conform that for $\gamma > 0$,

$$\int_{\mathcal{S}_{\mathbf{u}}^{\delta}} \left\{ \frac{p^*(\mathbf{s}|\mathbf{u})^{\gamma+1}}{\{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{s}|\mathbf{u})\}^{\gamma+1} + p(\mathbf{s})^{\gamma+1}} \right\}^{\frac{\gamma}{\gamma+1}} \delta(\mathbf{s}|\mathbf{u}) d\mathbf{s} = 0.$$

The above equation gives

$$\nu \leq O\left(\sup_{\mathbf{x}, \mathbf{u}} |r(\mathbf{x}, \mathbf{u}) - r^*(\mathbf{x}, \mathbf{u})|\right).$$

Thus, the proof is completed. \square

E Proof of Proposition 2

E.1 Main proof

Proof. We first define the following notations:

$$\begin{aligned}
L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) &= \frac{1}{1 + r_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{(\gamma+1)}} \\
S_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) &= \{L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})(1 - L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}))\}^{\gamma/(1+\gamma)} \\
\boldsymbol{\alpha}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) &= L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{1/(1+\gamma)} S_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) \frac{\partial \log r_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{\partial \boldsymbol{\theta}} \\
\boldsymbol{\beta}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) &= (1 - L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}))^{1/(1+\gamma)} S_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) \frac{\partial \log r_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{\partial \boldsymbol{\theta}} \\
\mathbf{C}_{\boldsymbol{\theta}} &= \iint \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \boldsymbol{\alpha}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{\top} \bar{\mathbf{p}}(\mathbf{x}, \mathbf{u}) - \frac{\partial}{\partial \boldsymbol{\theta}} \boldsymbol{\beta}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{\top} \bar{\mathbf{p}}(\mathbf{x}) \bar{\mathbf{p}}(\mathbf{u}) \right\} d\mathbf{x} d\mathbf{u}.
\end{aligned}$$

Then, our proof relies on the following lemma proved in Section E.2:

Lemma 1. *The influence function of $\hat{\boldsymbol{\theta}}$ under the γ -cross entropy (11) satisfies the following equations:*

- Under the contamination model (A),

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} \text{IF}(\bar{\mathbf{x}}) = \int \{ \boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \mathbf{u}) - \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \mathbf{u}) \} p^*(\mathbf{u}) d\mathbf{u}, \quad (43)$$

where IF denotes the influence function,

- Under the contamination model (B),

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} \text{IF}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) = -\boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) + \int \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}, \bar{\mathbf{u}}) p^*(\mathbf{x}) d\mathbf{x} + \int \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \mathbf{u}) p^*(\mathbf{u}) d\mathbf{u} - \iint \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}, \mathbf{u}) p^*(\mathbf{x}) p^*(\mathbf{u}) d\mathbf{x} d\mathbf{u}. \quad (44)$$

First of all, under Assumption (19), we observe that

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} \boldsymbol{\alpha}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) = \lim_{\|\mathbf{x}\| \rightarrow \infty} \boldsymbol{\beta}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u}) = \mathbf{0} \quad (45)$$

because $L_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})^{1/(1+\gamma)} \leq 1$. From Lemma 1, regarding the contamination model (A), (45) ensures that the right-hand side of (43) approaches to $\mathbf{0}$ as $\|\bar{\mathbf{x}}\| \rightarrow \infty$. Thus,

$$\lim_{\|\bar{\mathbf{x}}\| \rightarrow \infty} \|\text{IF}(\bar{\mathbf{x}})\| = \mathbf{0}.$$

Furthermore, under Assumption (18),

$$\sup_{\mathbf{x}, \mathbf{u}} \|\boldsymbol{\alpha}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})\| < \infty, \quad \sup_{\mathbf{x}, \mathbf{u}} \|\boldsymbol{\beta}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{u})\| < \infty. \quad (46)$$

This indicates that

$$\sup_{\bar{\mathbf{x}}} \|\text{IF}(\bar{\mathbf{x}})\| < \infty,$$

and therefore $\hat{\boldsymbol{\theta}}$ is B-robust.

For the the contamination model (B), (46) holds, which indicates that the right-hand side of (44) is uniformly bounded even when $r_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \bar{\mathbf{u}})$, $r_{\hat{\boldsymbol{\theta}}}(\bar{\mathbf{x}}, \mathbf{u})$ or $r_{\hat{\boldsymbol{\theta}}}(\mathbf{x}, \bar{\mathbf{u}})$ diverge through the influence of the outliers $(\bar{\mathbf{x}}, \bar{\mathbf{u}})$. Thus, we have

$$\sup_{\bar{\mathbf{x}}, \bar{\mathbf{u}}} \|\text{IF}(\bar{\mathbf{x}}, \bar{\mathbf{u}})\| < \infty. \quad (47)$$

□

E.2 Proof of Lemma 1

Proof. $\hat{\theta}$ associated with the densities $p^*(\mathbf{x}, \mathbf{u})$ and $p^*(\mathbf{x})p^*(\mathbf{u})$ is a solution of the following estimating equation:

$$\begin{aligned} \mathbf{0} &= \left. \frac{\partial}{\partial \boldsymbol{\theta}} d_\gamma(p^*(y|\mathbf{x}, \mathbf{u}), r_\theta(\mathbf{x}, \mathbf{u}); p^*(\mathbf{x}, \mathbf{u})) \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ &\propto \left[\iint \left\{ \left(\frac{r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}} \right)^{-1/(1+\gamma)} \frac{\frac{\partial}{\partial \boldsymbol{\theta}} r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{(1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)})^2} p^*(\mathbf{x}, \mathbf{u}) \right. \right. \\ &\quad \left. \left. - \left(\frac{1}{1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}} \right)^{-1/(1+\gamma)} \frac{\frac{\partial}{\partial \boldsymbol{\theta}} r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{(1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)})^2} p^*(\mathbf{x})p^*(\mathbf{u}) \right\} d\mathbf{x}d\mathbf{u} \right] \Bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ &= \iint \left\{ \boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x}, \mathbf{u}) - \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u}) \right\} d\mathbf{x}d\mathbf{u} \end{aligned} \quad (48)$$

On the other hand, $\hat{\boldsymbol{\theta}}_\epsilon$ is a solution of the the estimating equation with the contaminated densities:

$$\begin{aligned} \mathbf{0} &= \left. \frac{\partial}{\partial \boldsymbol{\theta}} d_\gamma(\bar{p}(y|\mathbf{x}, \mathbf{u}), r_\theta(\mathbf{x}, \mathbf{u}); \bar{p}(\mathbf{x}, \mathbf{u})) \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_\epsilon} \\ &\propto \left[\iint \left\{ \left(\frac{r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}} \right)^{-1/(1+\gamma)} \frac{\frac{\partial}{\partial \boldsymbol{\theta}} r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{(1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)})^2} \bar{p}(\mathbf{x}, \mathbf{u}) \right. \right. \\ &\quad \left. \left. - \left(\frac{1}{1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}} \right)^{-1/(1+\gamma)} \frac{\frac{\partial}{\partial \boldsymbol{\theta}} r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)}}{(1+r_\theta(\mathbf{x}, \mathbf{u})^{(\gamma+1)})^2} \bar{p}(\mathbf{x})\bar{p}(\mathbf{u}) \right\} d\mathbf{x}d\mathbf{u} \right] \Bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_\epsilon} \\ &= \iint \left\{ \boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}_\epsilon}(\mathbf{x}, \mathbf{u})\bar{p}(\mathbf{x}, \mathbf{u}) - \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}_\epsilon}(\mathbf{x}, \mathbf{u})\bar{p}(\mathbf{x})\bar{p}(\mathbf{u}) \right\} d\mathbf{x}d\mathbf{u}. \end{aligned} \quad (49)$$

Next, we obtain (43) and (44) under two contamination models separately by combining (49) with (48).

E.2.1 Contamination model (A)

We remind that $\bar{p}(\mathbf{x}, \mathbf{u})$, $\bar{p}(\mathbf{x})$, and $\bar{p}(\mathbf{u})$ in the contamination model (A) are defined as

$$\begin{aligned} \bar{p}(\mathbf{x}, \mathbf{u}) &= (1-\epsilon)p^*(\mathbf{x}, \mathbf{u}) + \epsilon\bar{\delta}_{\bar{\mathbf{x}}}(\mathbf{x})p^*(\mathbf{u}) \\ \bar{p}(\mathbf{x}) &= (1-\epsilon)p^*(\mathbf{x}) + \epsilon\bar{\delta}_{\bar{\mathbf{x}}}(\mathbf{x}) \\ \bar{p}(\mathbf{u}) &= p^*(\mathbf{u}), \end{aligned}$$

with $\bar{\delta}_{\bar{\mathbf{z}}}(z)$ is the Dirac delta function having a point mass at $\bar{\mathbf{z}}$. Then, the right-hand side of (49) is given by

$$(1-\epsilon) \underbrace{\iint \left\{ \boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x}, \mathbf{u}) - \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u}) \right\} d\mathbf{x}d\mathbf{u}}_{(*)} + \epsilon \int \left\{ \boldsymbol{\alpha}_{\hat{\boldsymbol{\theta}}_\epsilon}(\bar{\mathbf{x}}, \mathbf{u}) - \boldsymbol{\beta}_{\hat{\boldsymbol{\theta}}_\epsilon}(\bar{\mathbf{x}}, \mathbf{u}) \right\} p^*(\mathbf{u}) d\mathbf{u}. \quad (50)$$

Taylor expansion of (50) around $\hat{\boldsymbol{\theta}}$ is given by

$$\mathbf{0} = \left. \frac{\partial d_\gamma}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} + \left. \frac{\partial d_\gamma}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} (\hat{\boldsymbol{\theta}}_\epsilon - \hat{\boldsymbol{\theta}}) + O(\|\hat{\boldsymbol{\theta}}_\epsilon - \hat{\boldsymbol{\theta}}\|^2). \quad (51)$$

Since it follows from (48) that the term $(*)$ in (50) vanishes as $\hat{\boldsymbol{\theta}}_\epsilon = \hat{\boldsymbol{\theta}}$, we obtain (43) by taking the limit of $\epsilon \rightarrow 0$ because $\mathbf{C}_{\hat{\boldsymbol{\theta}}} = \left. \frac{\partial d_\gamma}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$ in (51).

E.2.2 Contamination model (B)

In the contamination model (B), the auxiliary variables \mathbf{u} are also contaminated as follows:

$$\begin{aligned}\bar{p}(\mathbf{x}, \mathbf{u}) &= (1 - \epsilon)p^*(\mathbf{x}, \mathbf{u}) + \epsilon\bar{\delta}_{\mathbf{x}}(\mathbf{x})\bar{\delta}_{\mathbf{u}}(\mathbf{u}) \\ \bar{p}(\mathbf{x}) &= (1 - \epsilon)p^*(\mathbf{x}) + \epsilon\bar{\delta}_{\mathbf{x}}(\mathbf{x}) \\ \bar{p}(\mathbf{u}) &= (1 - \epsilon)p^*(\mathbf{u}) + \epsilon\bar{\delta}_{\mathbf{u}}(\mathbf{u}).\end{aligned}$$

Then, the right-hand side of (49) is given by

$$\begin{aligned}& \iint \alpha_{\hat{\theta}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x}, \mathbf{u}) - \iint \beta_{\hat{\theta}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u})d\mathbf{x}d\mathbf{u} + \epsilon \left[\alpha_{\hat{\theta}_\epsilon}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) - \int \beta_{\hat{\theta}_\epsilon}(\mathbf{x}, \bar{\mathbf{u}})p^*(\mathbf{x})d\mathbf{x} \right. \\ & \left. - \int \beta_{\hat{\theta}_\epsilon}(\bar{\mathbf{x}}, \mathbf{u})p^*(\mathbf{u})d\mathbf{u} - \iint \alpha_{\hat{\theta}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x}, \mathbf{u}) + 2 \iint \beta_{\hat{\theta}_\epsilon}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u})d\mathbf{x}d\mathbf{u} \right] + O(\epsilon^2),\end{aligned}\quad (52)$$

By following the derivation in Section E.2.1, we have

$$\begin{aligned}C_{\hat{\theta}}\text{IF}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) &= -\alpha_{\hat{\theta}}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) + \int \beta_{\hat{\theta}}(\mathbf{x}, \bar{\mathbf{u}})p^*(\mathbf{x})d\mathbf{x} + \int \beta_{\hat{\theta}}(\bar{\mathbf{x}}, \mathbf{u})p^*(\mathbf{u})d\mathbf{u} \\ & \quad + \iint \alpha_{\hat{\theta}}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x}, \mathbf{u}) - 2 \iint \beta_{\hat{\theta}}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u})d\mathbf{x}d\mathbf{u} \\ &= -\alpha_{\hat{\theta}}(\bar{\mathbf{x}}, \bar{\mathbf{u}}) + \int \beta_{\hat{\theta}}(\mathbf{x}, \bar{\mathbf{u}})p^*(\mathbf{x})d\mathbf{x} + \int \beta_{\hat{\theta}}(\bar{\mathbf{x}}, \mathbf{u})p^*(\mathbf{u})d\mathbf{u} - \iint \beta_{\hat{\theta}}(\mathbf{x}, \mathbf{u})p^*(\mathbf{x})p^*(\mathbf{u})d\mathbf{x}d\mathbf{u},\end{aligned}$$

where (48) was applied in the last line. Thus, the proof is completed. \square

F Robust property in multiclass classification

Here, we perform a similar robust analysis as Theorem 2 in multiclass classification. We first specifically assume that the following $\tilde{\nu}$ is sufficiently small:

$$\tilde{\nu} := \int \frac{\sum_{\mathbf{u}} r(\mathbf{u}, \mathbf{x})^\gamma \delta(\mathbf{x}|\mathbf{u})\epsilon(\mathbf{u})p(\mathbf{u})}{(\sum_{\mathbf{u}'} r(\mathbf{u}', \mathbf{x})^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{x}.\quad (53)$$

Then, the outlier model (5) decomposes the γ -cross entropy (20) into

$$\begin{aligned}& d_\gamma(p(\mathbf{u}|\mathbf{x}), r(\mathbf{u}, \mathbf{x}); p(\mathbf{x})) \\ &= -\frac{1}{\gamma} \log \int \frac{\sum_{\mathbf{u}} r(\mathbf{u}, \mathbf{x})^\gamma p(\mathbf{u}|\mathbf{x})}{(\sum_{\mathbf{u}'} r(\mathbf{u}', \mathbf{x})^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} p(\mathbf{x})d\mathbf{x} \\ &= -\frac{1}{\gamma} \log \left[\int \frac{\sum_{\mathbf{u}} r(\mathbf{u}, \mathbf{x})^\gamma (1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})p(\mathbf{u})}{(\sum_{\mathbf{u}'} r(\mathbf{u}', \mathbf{x})^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{x} + \int \frac{\sum_{\mathbf{u}} r(\mathbf{u}, \mathbf{x})^\gamma \delta(\mathbf{x}|\mathbf{u})\epsilon(\mathbf{u})p(\mathbf{u})}{(\sum_{\mathbf{u}'} r(\mathbf{u}', \mathbf{x})^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{x} \right] \\ &= d_\gamma(r(\mathbf{u}, \mathbf{x}), p^*(\mathbf{x}|\mathbf{u}); (1 - \epsilon(\mathbf{u}))p(\mathbf{u})) + O(\tilde{\nu}),\end{aligned}\quad (54)$$

where we employed $\log(y + z) = \log(y) + O(z)$ with sufficiently small z and $p(\mathbf{u}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{x})} = \frac{(1 - \epsilon(\mathbf{u}))p^*(\mathbf{x}|\mathbf{u})p(\mathbf{u}) + \epsilon(\mathbf{u})\delta(\mathbf{x}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{x})}$, and

$$d_\gamma(r(\mathbf{u}, \mathbf{x}), p^*(\mathbf{x}|\mathbf{u}); (1 - \epsilon(\mathbf{u}))p(\mathbf{u})) := -\frac{1}{\gamma} \log \left[\sum_{\mathbf{u}} \left\{ \int \frac{r(\mathbf{u}, \mathbf{x})^\gamma p^*(\mathbf{x}|\mathbf{u})}{(\sum_{\mathbf{u}'} r(\mathbf{u}', \mathbf{x})^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{x} \right\} (1 - \epsilon(\mathbf{u}))p(\mathbf{u}) \right].$$

Eq.(54) indicates that minimization of $d_\gamma(p(\mathbf{u}|\mathbf{x}), r(\mathbf{u}, \mathbf{x}); p(\mathbf{x}))$ approximately equals to minimization of $d_\gamma(r(\mathbf{u}, \mathbf{x}), p^*(\mathbf{x}|\mathbf{u}); (1 - \epsilon(\mathbf{u}))p(\mathbf{u}))$ when $\tilde{\nu}$ is sufficiently small. In addition, $d_\gamma(r(\mathbf{u}, \mathbf{x}), p^*(\mathbf{x}|\mathbf{u}); (1 - \epsilon(\mathbf{u}))p(\mathbf{u}))$ is minimized at

$$r(\mathbf{u}, \mathbf{x}) = p^*(\mathbf{x}|\mathbf{u})\quad (55)$$

because it is the γ -cross entropy to $p^*(\mathbf{x}|u)$ under the measure $(1 - \epsilon(u))p(u)$. The minimizer (55) is desirable in terms of the universal approximation assumptions (A4) and (B4).

Next, we discuss when $\tilde{\nu}$ is sufficiently small. Following Section D, in the neighborhood of $p^*(\mathbf{x}|u)$,

$$\begin{aligned}\tilde{\nu} &\leq \sum_u \left\{ \int \frac{p^*(\mathbf{x}|u)^\gamma \delta(\mathbf{x}|u)}{(\sum_{u'} p^*(\mathbf{x}|u')^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{x} \right\} \epsilon(u)p(u) + O(\sup_{u,\mathbf{x}} |r(u,\mathbf{x}) - p^*(\mathbf{x}|u)|) \\ &= \sum_u \left\{ \int \frac{p^*(\mathbf{s}|u)^\gamma \delta(\mathbf{s}|u)}{(\sum_{u'} p^*(\mathbf{s}|u')^{\gamma+1})^{\frac{\gamma}{\gamma+1}}} d\mathbf{s} \right\} \epsilon(u)p(u) + O(\sup_{u,\mathbf{x}} |r(u,\mathbf{x}) - p^*(\mathbf{x}|u)|),\end{aligned}\quad (56)$$

where we performed the change of variables from \mathbf{x} to \mathbf{s} under the data generate model (1). When $\sum_{u'} p^*(\mathbf{s}|u')^{\gamma+1} \neq 0$ and the supports of $p^*(\mathbf{s}|u)$ and $\delta(\mathbf{s}|u)$ are mutually disjoint as in Section D, we can make the same implication as Proposition 1: $\tilde{\nu}$ can be sufficiently small in the neighborhood of $p^*(\mathbf{x}|u)$ when $p^*(\mathbf{s}|u)$ and $\delta(\mathbf{s}|u)$ are clearly separated. This clear separation possibly happens on a situation where $\delta(\mathbf{s}|u)$ lies on the tails of $p^*(\mathbf{s}|u)$ as seen in common contamination by outliers. Thus, the γ -cross entropy for multiclass classification would be also robust against outliers.

As discussed in Section 4.3, the non-robustness of TCL can be understood in terms of the γ -cross entropy. TCL employs the multiclass logistic regression whose cross-entropy can be obtained as a limit of $\gamma = 0$ in the γ -cross entropy (20). The multiclass logistic regression cannot also fulfill the robustness condition: It follows from the definition (53) that $\tilde{\nu}$ cannot be sufficiently small when $\gamma = 0$. Thus, this implies that TCL is sensitive to outliers.

G Experimental details

G.1 Robust time contrastive learning

Source vectors with time segment length 512 was first generated from (5): Following (2), given a time segment label u , the target density $p^*(\mathbf{s}|u)$ was conditionally independent Laplace densities with means 0 and different scales across time segments, which were randomly determined from the uniform density on $[0, \frac{1}{\sqrt{2}}]$. Regarding the outlier density $\delta(\mathbf{s}|u)$, two types of densities were used: An independent Laplace density with mean 0 and scale 3.0, and a mixture of mean-modulated two Gaussians. More precisely, the outlier density of the mixture is given by $\delta(\mathbf{s}|u) = 0.5N(1+3y(u), 0.5) + 0.5N(-1-3y(u), 0.5)$ where $N(a, b)$ is a Gaussian density with mean a and standard deviation b , and $y(u)$ is fixed at the scale parameter of $p^*(\mathbf{s}|u)$ above in the u -th time segment. We set $\epsilon(u) = \epsilon$ for all time segments u . The total numbers of segments and of data samples were $K = 256$ and $T = 512 \times 256$, respectively. The dimensionality of data was $d_x = 10$ in Table 1, while $d_x = 5$ in Table 2. Then, data \mathbf{x} was generated according to (1) where $\mathbf{f}(\mathbf{s})$ was modelled by a three-layer neural network (Table 1) or two-layer (Table 2) neural network with the leaky ReLU activation function and random weights. The numbers of all hidden and output units were the same as the dimensionality of data (i.e., d_x). As preprocessing, we performed whitening based on the γ -divergence [Chen et al., 2013].

ICA features $\mathbf{h}(\mathbf{x})$ both in RTCL and TCL were modelled by a three layer neural network where the number of hidden units was $4d_x$, but the final layer was d_x . Regarding the activation functions, the final layer employed the absolute value function, while the other hidden layers were the max-out function [Goodfellow et al., 2013] with two groups. ℓ_2 regularization was employed with the regularization parameter 10^{-4} . When $\epsilon < 0.1$, we optimized the network parameters both in RTCL and TCL with 0.001 learning rate using the Adam optimizer for 1,000 epochs with mini-batch size 256, while we updated the parameters for 3,000 epochs for $\epsilon = 0.1$. We empirically observed that more iterations for parameter update are often needed to escape from bad local optima when the contamination ratio is larger. No postprocessing was applied. The performance was measured by the absolute value of the Pearson correlation coefficient between learned ICA features $\mathbf{h}(\mathbf{x})$ and \mathbf{s} without outliers.

When comparing with iVAE, the same neural network has been employed both in RTCL and iVAE, but the activation functions were all ReLU except for the final layer. We optimized the network parameters with 0.001 learning rate using the Adam optimizer for 1,000 epochs, and the mini-batch size was 256.

G.2 Robust permutation contrastive learning

First, the temporally dependent T sources were generated from

$$\log p^*(\mathbf{s}(t)|\mathbf{s}(t-1)) = -\sum_{i=1}^{d_x} |s_i(t) - \rho s_i(t-1)| + C,$$

where C denotes a constant and the auto-regressive coefficient ρ was fixed at 0.7. The total number of sources was $T = 65,536$. Then, we randomly replaced the sources with outliers based on a constant contamination ratio ϵ , which were generated from the independent Gaussian density with the mean 0 and scale 3. Data \mathbf{x} was generated as the nonlinear mixing of the sources with outliers by a three-layer neural networks with the leaky ReLU function and random weights according to (1). The same whitening preprocessing above was performed based on the γ -divergence [Chen et al., 2013].

Based on the universal approximation assumption in Hyvärinen and Morioka [2017, Theorem 1 and Eq.(12)] or Appendix B, we restrict the model r as

$$r(\mathbf{x}(t), \mathbf{u}(t)) = \exp\left(\sum_{i=1}^{d_x} \psi_i(h_i(\mathbf{x}(t)), h_i(\mathbf{u}(t)))\right),$$

with a neural network $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_{d_x}(\mathbf{x}))^\top$. Following Hyvärinen and Morioka [2017], $\psi_i(h_i(\mathbf{x}), h_i(\mathbf{u}))$ was further modelled by

$$|a_{i,1}h_i(\mathbf{x}) + a_{i,2}h_i(\mathbf{u}) + b_i| - (\bar{a}_i h_i(\mathbf{x}) + \bar{b}_i)^2 + c,$$

where $a_{i,1}, a_{i,2}, b_i, \bar{a}_i, \bar{b}_i, c$ are parameters to be estimated from data. ICA features $\mathbf{h}(\mathbf{x})$ both in RPCL and PCL were modelled by a three layer feedforward neural network where the number of hidden units is $4d_x$. No activation function was applied in the last layer, while the max-out function was employed in the other layers. ℓ_2 regularization was employed with the regularization parameter 10^{-4} . Then, we optimized the parameters in RPCL and PCL using the Adam optimizer with 0.001 learning rate for 1,000 epochs with mini-batch size 128. The performance was measured by the absolute Pearson correlation coefficient between learned ICA features $\mathbf{h}(\mathbf{x})$ and \mathbf{s} without outliers.

H Visualization of outliers on fMRI data

Fig.1 visualizes the time series data from the Parahippocampal brain region, and demonstrates the presence of outliers in many of the segments.

References

- P. Chen, H. Hung, O. Komori, S.-Y. Huang, and S. Eguchi. Robust independent component analysis via minimum γ -divergence estimation. *IEEE Journal of Selected Topics in Signal Processing*, 7(4):614–624, 2013.
- I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pages 1319–1327. PMLR, 2013.
- A. Hyvärinen and H. Morioka. Nonlinear ICA of temporally dependent stationary sources. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54, pages 460–469. PMLR, 2017.
- A. Hyvärinen, H. Sasaki, and R. E. Turner. Nonlinear ICA using auxiliary variables and generalized contrastive learning. In *Proceedings of the 22th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 89, pages 859–868, 2019.

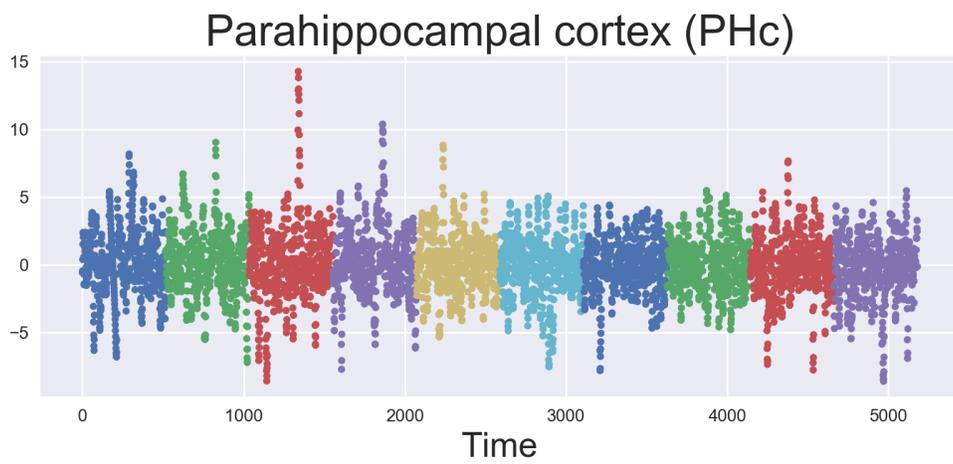


Figure 1: Subset of time series data corresponding to the Parahippocampal (PHc) region taken from the Hippocampal fMRI dataset. Different colors denote distinct segments, which in this case correspond to fMRI measurements from the same subject on distinct days.