

Density estimation based on U -divergence

Osamu Komori Prediction and Knowledge Discovery Research Center, Project Researcher

1 U -divergence

Let $U : \mathbf{R}^+ \rightarrow \mathbf{R}$ be a convex and strictly increasing function with the derivative u and the inverse function $\xi = u^{-1}$. Then for real-valued functions f and $g : \mathbf{R}^p \rightarrow \mathbf{R}^+$, the U -divergence is given as a special case of the Bregman divergence (?):

$$D_U(g, f) = \int d(\xi(g(\mathbf{x})), \xi(f(\mathbf{x})))d\mathbf{x}, \quad (1)$$

where

$$d(g', f') = U(f') - \{u(g')(f' - g') + U(g')\}. \quad (2)$$

Note that $D_U(g, f)$ is non-negative because of the convexity of U . The equality holds if and only if $f = g$ (a.e. \mathbf{x}). It is also simply expressed as

$$D_U(g, f) = C_U(g, f) - H_U(g), \quad (3)$$

where

$$C_U(g, f) = - \int g(\mathbf{x})\xi(f(\mathbf{x}))d\mathbf{x} + \int U(\xi(f(\mathbf{x})))d\mathbf{x} \quad (4)$$

$$H_U(g) = - \int g(\mathbf{x})\xi(g(\mathbf{x}))d\mathbf{x} + \int U(\xi(g(\mathbf{x})))d\mathbf{x} (= C_U(g, g)), \quad (5)$$

and $C_U(g, f)$ and $H_U(g)$ are called the U -cross entropy and U -entropy, respectively.

U -loss function with volume-mass-one

The U -loss function for observations $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, which derived from the cross-entropy in (4), is defined as

$$L_U(f) = -\frac{1}{n} \sum_{i=1}^n \xi(f(\mathbf{x}_i)) + \int U(\xi(f(\mathbf{x})))d\mathbf{x}. \quad (6)$$

Then, we consider the following variant:

$$\mathcal{L}_U(f) \equiv L_U(u(U^{-1}(f))) \quad (7)$$

$$= -\frac{1}{n} \sum_{i=1}^n U^{-1}(f(\mathbf{x}_i)) + 1 \quad (8)$$

The point is that the second integral term in (6) is restricted to be 1, which we call volume-mass-one. Here we consider $U(t) = (1 + \beta t)^{(1+\beta)/\beta} / (1 + \beta)$ with $\beta > 0$.

2 Algorithm

1. Set $f_0(\mathbf{x}) = 0$.

2. For $k = 1, \dots, K$,

(a) Initialize $\pi = \pi_0 (\ll 1)$, $\Sigma = \mathbf{I}$ and $\boldsymbol{\mu} = \operatorname{argmin}_{\boldsymbol{\mu} \in D} \left\{ \mathcal{L}_\beta \left((1 - \pi)f_{k-1}^{1+\beta} + \pi\phi(\boldsymbol{\mu}, \mathbf{I}) \right) \right\}$, where \mathbf{I} is the $p \times p$ identity matrix; ϕ is the basis function in \mathcal{D}_β . Define

$$\mathcal{R}_{\boldsymbol{\mu}, \Sigma} = \left\{ i \mid \frac{\beta}{2(1+\beta)} (\mathbf{x}_i - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) < 1, \mathbf{x}_i \in D \right\}. \quad (9)$$

(b) For \mathbf{x}_i such that $i \in \mathcal{R}_{\boldsymbol{\mu}, \Sigma}$, calculate

$$q(\mathbf{x}_i) = \frac{\pi\phi(\mathbf{x}_i)}{(1 - \pi)f_{k-1}^{1+\beta}(\mathbf{x}_i) + \pi\phi(\mathbf{x}_i)} \quad (10)$$

$$\boldsymbol{\mu}_q = \frac{\sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} q(\mathbf{x}_i)^{\frac{1}{1+\beta}} \mathbf{x}_i}{\sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} q(\mathbf{x}_i)^{\frac{1}{1+\beta}}}. \quad (11)$$

where $\sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}}$ is the summation of i over $\mathcal{R}_{\boldsymbol{\mu}, \Sigma}$.

(c) Update $\boldsymbol{\mu} = \boldsymbol{\mu}_q$ and go to step (d) if $\mathcal{R}_{\boldsymbol{\mu}, \Sigma} \subset \mathcal{R}_{\boldsymbol{\mu}_q, \Sigma}$; otherwise go back to step (b).

(d) For \mathbf{x}_i such that $i \in \mathcal{R}_{\boldsymbol{\mu}, \Sigma}$, update $q(\mathbf{x}_i)$ as in (10) and calculate

$$\Sigma_q = \frac{2 + (2 + p)\beta \sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} q(\mathbf{x}_i)^{\frac{1}{1+\beta}} (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})'}{2(1 + \beta) \sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} q(\mathbf{x}_i)^{\frac{1}{1+\beta}}}. \quad (12)$$

(e) Update $\Sigma = \Sigma_q$ and go to step (f) if $\mathcal{R}_{\boldsymbol{\mu}, \Sigma} \subset \mathcal{R}_{\boldsymbol{\mu}, \Sigma_q}$; otherwise go back to step (d).

(f) For \mathbf{x}_i such that $i \in \mathcal{R}_{\boldsymbol{\mu}, \Sigma}$, update $q(\mathbf{x}_i)$ as in (10) and calculate

$$\pi_q = \frac{A_2^{1+\beta}}{A_1^{1+\beta} + A_2^{1+\beta}}, \quad (13)$$

where

$$A_1 = \sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} (1 - q(\mathbf{x}_i))^{\frac{1}{1+\beta}} f_{k-1}(\mathbf{x}_i)^\beta \quad (14)$$

$$A_2 = \sum_{\mathcal{R}_{\boldsymbol{\mu}, \Sigma}} q(\mathbf{x}_i)^{\frac{1}{1+\beta}} \phi(\mathbf{x}_i)^{\frac{\beta}{1+\beta}}, \quad (15)$$

and update $\pi = \pi_q$, and $q(\mathbf{x}_i)$ as in (10).

(g) Repeat the steps from (b) to (f) until the values of $\boldsymbol{\mu}$, Σ and π converges, and set them to be $\boldsymbol{\mu}_k$, Σ_k , π_k , respectively.

(h) Update f_{k-1} with $\phi_k(\mathbf{x}) = \phi_\beta(\mathbf{x}, \boldsymbol{\mu}_k, \Sigma_k)$ and π_k as

$$f_k = \left\{ (1 - \pi_k) f_{k-1}^{1+\beta} + \pi_k \phi_k \right\}^{\frac{1}{1+\beta}}. \quad (16)$$

3. Output $\hat{f} = f_K$.

Theorem 2.1 The empirical loss $\mathcal{L}_\beta(f_k)$ in the boosting algorithm is monotonically decreasing with respect to k . That is, for $k = 1, \dots, K$,

$$\mathcal{L}_\beta(f_k) \leq \mathcal{L}_\beta(f_{k-1}). \quad (17)$$

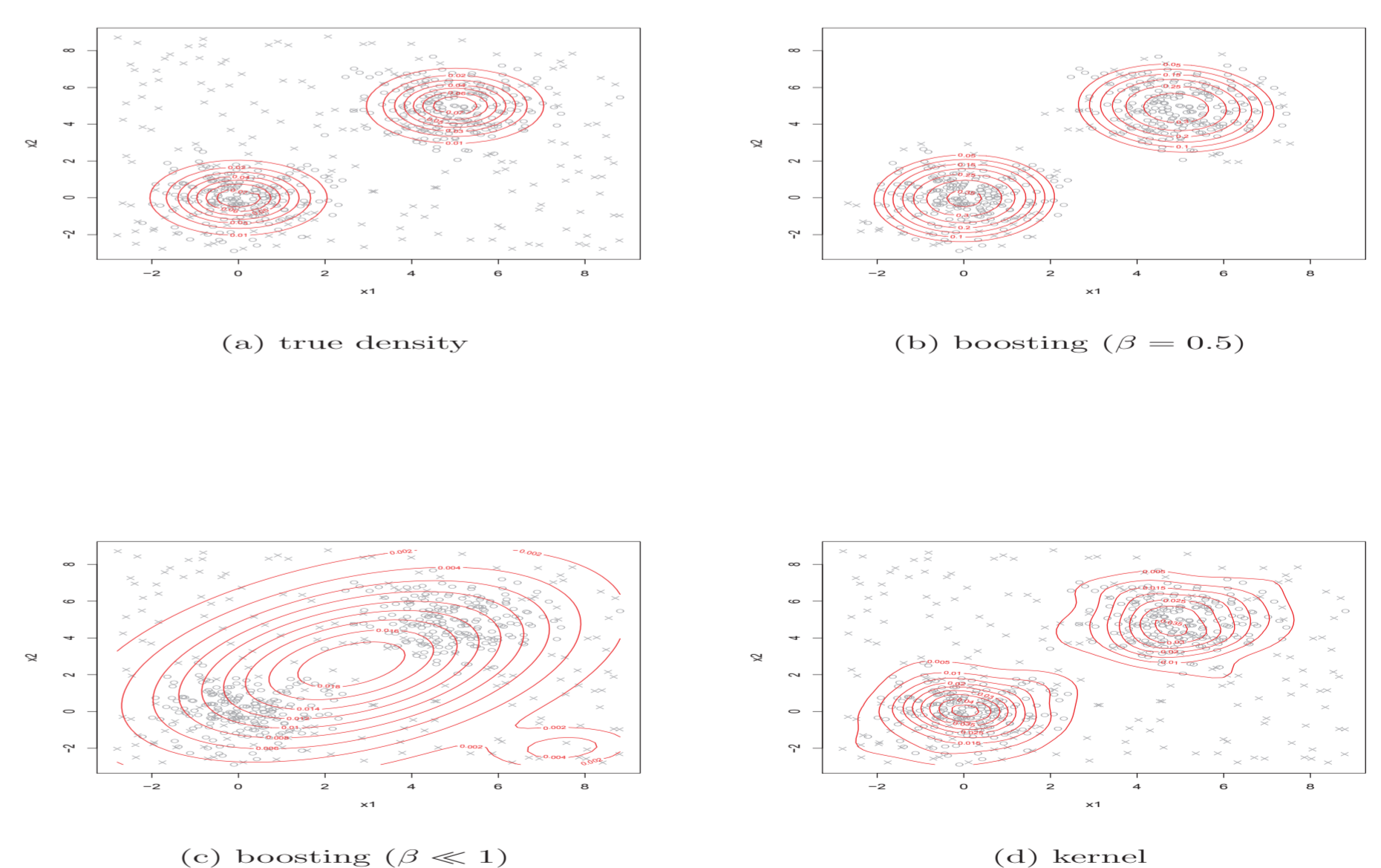


Fig1. Contour plots for the true density (a) and density estimators by three methods (b), (c) and (d). Observations from the normal distributions are denoted by circles; noisy observations are denoted by cross marks. Observations that are not used in the estimation are deleted in the panel (b).