

Supplementary Material of Robust Bayesian Kernel Machine via Stein Variational Gradient Descent for Big Data

January 26, 2018

1 Loss functions for multiclass problem

Given a data example (\mathbf{x}, y) and M hyperplanes $\mathbf{w}_1, \dots, \mathbf{w}_M$, each corresponds to a class label. For notation convenience, we stack these vectors into a matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_M]^\top$. The discriminative score for the data example \mathbf{x} given by each class is denoted as $g_m(\mathbf{x})$. In the linear classification, $g_m(\mathbf{x})$ is defined as $\mathbf{w}_m^\top \mathbf{x}$. To deal with nonlinear nature of data, the input data \mathbf{x} is mapped onto a feature space via a transformation $\phi(\mathbf{x})$, thus the score $g_m(\mathbf{x})$ is of the form $\mathbf{w}_m^\top \phi(\mathbf{x})$. There are two popular loss functions for multiclass problem: the multiclass Hinge loss and the cross-entropy loss with softmax function.

1.1 Multiclass Hinge loss

Inspired from the binary Hinge loss, the multiclass Hinge loss is defined over the gap between the scores given by the true class and the runner-up as follows:

$$l(\mathbf{W}; \mathbf{x}, y) = \max \left\{ 0, 1 + \max_{j \neq y} g_j(\mathbf{x}) - g_y(\mathbf{x}) \right\}$$

This is a convex and non-smooth loss function, hence we have a set of sub-gradients at each point and sub-gradients at different points are typically not smoothly varied.

1.2 Cross-entropy loss with softmax function

Let us consider the one-hot vector \mathbf{y} of M elements where the y -th element is one. The one-hot vector \mathbf{y} describes the true distribution p of labels for the data example \mathbf{x} . The distribution q of labels for \mathbf{x} given by the current model \mathbf{W} is defined over the softmax function as follows:

$$q(Y = j | \mathbf{W}, \mathbf{x}) = \frac{\exp(g_j(\mathbf{x}))}{\sum_{m=1}^M \exp(g_m(\mathbf{x}))}$$

The cross-entropy between two discrete distribution p and q is defined as follows:

$$H(p, q) = - \sum_Y p(Y) \log q(Y | \mathbf{W}, \mathbf{x})$$

It measures the distance between p and q , thus it can be used as a loss function. If p and q is exactly the same, the loss function is equal to zero. Because $p(Y)$ equals 1 only when $Y = y$ and otherwise is zero, we can derive as follows:

$$l(\mathbf{W}; \mathbf{x}, y) = H(p, q) = - \log q(Y = y | \mathbf{W}, \mathbf{x}) = - \log \frac{\exp(g_y(\mathbf{x}))}{\sum_{m=1}^M \exp(g_m(\mathbf{x}))}$$

2 Gradients of the log of the posterior distribution

Given the training set $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ where $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ and $y \in \mathcal{Y}$. In this section, we consider the nonlinear case where the data are mapped onto a random feature space via a transformation ϕ defined as:

$$\phi(\mathbf{x}) = \frac{1}{\sqrt{D}} [\cos(\boldsymbol{\omega}_i^\top \mathbf{x}), \sin(\boldsymbol{\omega}_i^\top \mathbf{x})]_{i=1}^D$$

in which the random frequencies $\omega_i, i = 1, \dots, D$ is computed as follows:

$$\omega_i = \gamma \odot \mathbf{e}_i$$

where \mathbf{e}_i is drawn from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. We note that it is straightforward to derive the formula for the linear case from the nonlinear case where we simply set $\phi(\mathbf{x}) = \mathbf{x}$.

The posterior distribution of the model \mathbf{W} and kernel parameter γ defined as follows:

$$p(\mathbf{W}, \gamma \mid \mathcal{D}) \propto p(\mathcal{D} \mid \mathbf{W}, \gamma) p(\mathbf{W}) p(\gamma)$$

where we have defined

$$p(\mathcal{D} \mid \mathbf{W}, \gamma) \propto \prod_{(\mathbf{x}, y) \in \mathcal{D}} \exp(-p(\mathbf{x}, y) l(\mathbf{W}; \mathbf{x}, y))$$

and $p(\mathbf{W})$, $p(\gamma)$ is prior distributions. The log of the posterior distribution is of the form:

$$\log p(\mathbf{W}, \gamma \mid \mathcal{D}) = \log p(\mathbf{W}) + \log p(\gamma) - \mathbb{E}_{P_{\mathcal{X} \times \mathcal{Y}}} [l(\mathbf{W}; \mathbf{x}, y)] + \text{const}$$

We can approximate $\mathbb{E}_{P_{\mathcal{X} \times \mathcal{Y}}} [l(\mathbf{W}; \mathbf{x}, y)]$ by sampling a mini-batch \mathcal{B} from the training set \mathcal{D} as follows:

$$\mathbb{E}_{P_{\mathcal{X} \times \mathcal{Y}}} [l(\mathbf{W}; \mathbf{x}, y)] \approx \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} l(\mathbf{W}; \mathbf{x}, y)$$

In summary, the log of posterior distribution of model and kernel parameters is of the form:

$$\log p(\mathbf{W}, \gamma \mid \mathcal{D}) \approx \log p(\mathbf{W}, \gamma \mid \mathcal{B}) = \log p(\mathbf{W}) + \log p(\gamma) - \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} l(\mathbf{W}; \mathbf{x}, y) + \text{const}$$

In this section, our aim is to compute the gradients of $\log p(\mathbf{W}, \gamma \mid \mathcal{D})$ w.r.t to the model \mathbf{W} and kernel parameter γ respectively. In our current framework, the prior of \mathbf{W} is a Multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbf{I})$ and the prior of γ is Truncate Multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I}_d, \mathbf{0}, +\infty)$. Obviously, the gradients is of the form:

$$\begin{aligned} \nabla_{\mathbf{W}} \log p(\mathbf{W}, \gamma \mid \mathcal{B}) &= -\lambda \|\mathbf{W}\|_{2,2}^2 - \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \nabla_{\mathbf{W}} l(\mathbf{W}; \mathbf{x}, y) \\ \nabla_{\gamma} \log p(\mathbf{W}, \gamma \mid \mathcal{B}) &= -\frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \nabla_{\gamma} l(\mathbf{W}; \mathbf{x}, y) \end{aligned}$$

Next, we present the gradients of the loss function w.r.t the model \mathbf{W} and kernel parameter γ respectively.

• Compute $\nabla_{\mathbf{W}} l(\mathbf{W}; \mathbf{x}, y)$

For the Hinge loss function. We have $\nabla_{\mathbf{W}} l(\mathbf{W}; \mathbf{x}, y) = 0$ if $l(\mathbf{W}; \mathbf{x}, y) = 0$. Otherwise, we have

$$\nabla_{\mathbf{W}_m} l(\mathbf{W}; \mathbf{x}, y) = \begin{cases} -\phi(\mathbf{x}) & \text{if } m = y \\ \phi(\mathbf{x}) & \text{if } m = m^* \\ \mathbf{0} & \text{otherwise} \end{cases}$$

where $m^* = \underset{m \in \mathcal{Y} \setminus y}{\operatorname{argmax}} g_m(\mathbf{x}) = \underset{m \in \mathcal{Y} \setminus y}{\operatorname{argmax}} \mathbf{w}_m^\top \mathbf{x}$.

For cross-entropy with softmax function. We have

$$l(\mathbf{W}; \mathbf{x}, y) = -\log \frac{\exp(g_y(\mathbf{x}))}{\sum_{m=1}^M \exp(g_m(\mathbf{x}))} = -\log \frac{\exp(\mathbf{w}_y^\top \phi(\mathbf{x}))}{\sum_{m=1}^M \exp(\mathbf{w}_m^\top \phi(\mathbf{x}))} = -\mathbf{w}_y^\top \phi(\mathbf{x}) + \log \left(\sum_{m=1}^M \exp(\mathbf{w}_m^\top \phi(\mathbf{x})) \right)$$

The gradient of the loss function w.r.t \mathbf{W} is formulated as follows:

$$\begin{aligned} \nabla_{\mathbf{W}_y} l(\mathbf{W}; \mathbf{x}, y) &= -\phi(\mathbf{x}) + \frac{\phi(\mathbf{x}) \exp(\mathbf{w}_y^\top \phi(\mathbf{x}))}{\sum_{m=1}^M \exp(\mathbf{w}_m^\top \phi(\mathbf{x}))} = \phi(\mathbf{x}) \left(s(\mathbf{W}^\top \phi(\mathbf{x}))_y - 1 \right) \\ \nabla_{\mathbf{W}_{i \neq y}} l(\mathbf{W}; \mathbf{x}, y) &= \frac{\phi(\mathbf{x}) \exp(\mathbf{w}_i^\top \phi(\mathbf{x}))}{\sum_{m=1}^M \exp(\mathbf{w}_m^\top \phi(\mathbf{x}))} = \phi(\mathbf{x}) s(\mathbf{W}^\top \phi(\mathbf{x}))_i \end{aligned}$$

where $s(\mathbf{g})_i = \frac{\exp(g_i)}{\sum_{m=1}^M \exp(g_m)}$ and $\mathbf{g} = [g_1, \dots, g_M]$.

• **Compute $\nabla_{\gamma} l(\mathbf{W}; \mathbf{x}, y)$**

The gradient of the loss function w.r.t the kernel parameter γ is formulated as follows:

$$\nabla_{\gamma} l(\mathbf{W}; \mathbf{x}, y) = \nabla_{\gamma} \boldsymbol{\omega} \times \nabla_{\boldsymbol{\omega}} \phi(\mathbf{x}) \times \nabla_{\phi} l(\mathbf{W}; \mathbf{x}, y)$$

where we have

$$\begin{aligned} \nabla_{\gamma} \boldsymbol{\omega} &= [\nabla_{\gamma} \boldsymbol{\omega}_1, \dots, \nabla_{\gamma} \boldsymbol{\omega}_D] \\ \nabla_{\gamma} \boldsymbol{\omega}_i &= \mathbf{e}_i \\ \nabla_{\boldsymbol{\omega}} \phi(\mathbf{x}) &= [\nabla_{\boldsymbol{\omega}_1} \phi(\mathbf{x}), \dots, \nabla_{\boldsymbol{\omega}_D} \phi(\mathbf{x})]^T \\ \nabla_{\boldsymbol{\omega}_i} \phi(\mathbf{x}) &= [\mathbf{0}, -\mathbf{x} \sin(\boldsymbol{\omega}_i^T \mathbf{x}), \mathbf{x} \cos(\boldsymbol{\omega}_i^T \mathbf{x}), \mathbf{0}] \end{aligned}$$

Thus, we have

$$\nabla_{\gamma} l(\mathbf{W}; \mathbf{x}, y) = [-\mathbf{e}_i \odot \mathbf{x} \sin(\boldsymbol{\omega}_i^T \mathbf{x}), \mathbf{e}_i \odot \mathbf{x} \cos(\boldsymbol{\omega}_i^T \mathbf{x})]_{i=1}^D \times \nabla_{\phi} l(\mathbf{W}; \mathbf{x}, y)$$

For the Hinge loss function. We have

$$\begin{aligned} \nabla_{\phi} l(\mathbf{W}; \mathbf{x}, y) &= \nabla_{\phi} \left(\max \left\{ 0, 1 + \max_{m \neq y} \mathbf{w}_m^T \phi(\mathbf{x}) - \mathbf{w}_y^T \phi(\mathbf{x}) \right\} \right) \\ &= \mathbb{I}_{l(\mathbf{W}; \mathbf{x}, y) > 0} (\mathbf{w}_{m^*} - \mathbf{w}_y) \end{aligned}$$

For cross-entropy with softmax function. We have

$$\begin{aligned} \nabla_{\phi} l(\mathbf{W}; \mathbf{x}, y) &= \nabla_{\phi} \left(-\mathbf{w}_y^T \phi(\mathbf{x}) + \log \left(\sum_{m=1}^M \exp(\mathbf{w}_m^T \phi(\mathbf{x})) \right) \right) \\ &= -\mathbf{w}_y + \frac{\sum_{m=1}^M \mathbf{w}_m \exp(\mathbf{w}_m^T \phi(\mathbf{x}))}{\sum_{m=1}^M \exp(\mathbf{w}_m^T \phi(\mathbf{x}))} \\ &= -\mathbf{w}_y + \sum_{m=1}^M \mathbf{w}_m s(\mathbf{W}^T \phi(\mathbf{x}))_m \end{aligned}$$