

## Supplementary Materials

### Localized Lasso for High-Dimensional Regression

Makoto Yamada<sup>1,2</sup>, Koh Takeuchi<sup>3</sup>, Tomoharu Iwata<sup>3</sup>, John Shawe-Taylor<sup>4</sup>, Samuel Kaski<sup>5</sup>

<sup>1</sup>RIKEN AIP, Japan

<sup>2</sup>JST, PRESTO

<sup>3</sup>NTT Communication Science Laboratories, Japan

<sup>4</sup>University College London, UK

<sup>5</sup>Aalto University, Finland

`makoto.yamada@riken.jp, {takeuchi.koh, iwata.tomoharu}@lab.ntt.co.jp  
j.shawe-taylor@ucl.ac.uk, samuel.kaski@aalto.fi`

### Propositions used for deriving Eq. (4) in main paper

**Proposition 1** Under  $r_{ij} \geq 0$ ,  $r_{ij} = r_{ji}$ ,  $r_{ii} = 0$ , we have

$$\frac{\partial}{\partial \text{vec}(\mathbf{W})} \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i - \mathbf{w}_j\|_2 = 2\mathbf{F}_g \text{vec}(\mathbf{W}),$$

where

$$\begin{aligned} \mathbf{F}_g &= \mathbf{I}_d \otimes \mathbf{C}, \\ [\mathbf{C}]_{i,j} &= \begin{cases} \frac{\sum_{j'=1}^n \frac{r_{ij'}}{\|\mathbf{w}_i - \mathbf{w}'_j\|_2} - \frac{r_{ij}}{\|\mathbf{w}_i - \mathbf{w}_j\|_2}}{\|\mathbf{w}_i - \mathbf{w}_j\|_2} & (i = j) \\ \frac{-r_{ij}}{\|\mathbf{w}_i - \mathbf{w}_j\|_2} & (i \neq j) \end{cases}. \end{aligned}$$

*Proof:* Under  $r_{ij} \geq 0$ ,  $r_{ij} = r_{ji}$ ,  $r_{ii} = 0$ , the derivative of the network regularization term with respect to  $\mathbf{w}_k$  is given as

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}_k} \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i - \mathbf{w}_j\|_2 &= \sum_{i=1}^n r_{ik} \frac{\mathbf{w}_k - \mathbf{w}_i}{\|\mathbf{w}_k - \mathbf{w}_i\|_2} + \sum_{j=1}^n r_{kj} \frac{\mathbf{w}_k - \mathbf{w}_j}{\|\mathbf{w}_j - \mathbf{w}_k\|_2} \\ &= \mathbf{w}_k \left( \sum_{i=1}^n \frac{r_{ik}}{\|\mathbf{w}_k - \mathbf{w}_i\|_2} + \sum_{j=1}^n \frac{r_{kj}}{\|\mathbf{w}_j - \mathbf{w}_k\|_2} \right) \\ &\quad - \sum_{i=1}^n \frac{r_{ik}}{\|\mathbf{w}_k - \mathbf{w}_i\|_2} \mathbf{w}_i - \sum_{j=1}^n \frac{r_{kj}}{\|\mathbf{w}_j - \mathbf{w}_k\|_2} \mathbf{w}_j \\ &= 2 \left( \mathbf{w}_k \sum_{i=1}^n \frac{r_{ik}}{\|\mathbf{w}_k - \mathbf{w}_i\|_2} - \sum_{i=1}^n \frac{r_{ik}}{\|\mathbf{w}_k - \mathbf{w}_i\|_2} \mathbf{w}_i \right). \end{aligned}$$

Thus,

$$\frac{\partial}{\partial \mathbf{W}} \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i - \mathbf{w}_j\|_2 = 2\mathbf{C}\mathbf{W},$$

where  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_n]^\top \in \mathbb{R}^{n \times d}$ . Since  $\text{vec}(\mathbf{C}\mathbf{W}\mathbf{I}_d) = (\mathbf{I}_d \otimes \mathbf{C})\text{vec}(\mathbf{W})$ , we have

$$\frac{\partial}{\partial \text{vec}(\mathbf{W})} \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i - \mathbf{w}_j\|_2 = 2(\mathbf{I}_d \otimes \mathbf{C})\text{vec}(\mathbf{W}),$$

where  $\mathbf{I}_d \in \mathbb{R}^{d \times d}$  is the identity matrix and  $\text{vec}(\cdot)$  is the vectorization operator.  $\square$

### Proposition 2

$$\frac{\partial}{\partial \text{vec}(\mathbf{W})} \sum_{i=1}^n \|\mathbf{w}_i\|_1^2 = 2\mathbf{F}_e \text{vec}(\mathbf{W}),$$

where

$$[\mathbf{F}_e]_{\ell,\ell} = \sum_{i=1}^n \frac{I_{i,\ell} \|\mathbf{w}_i\|_1}{[\text{vec}(|\mathbf{W}|)]_\ell}.$$

Hence,  $I_{i,\ell} \in \{0, 1\}$  are the group index indicators:  $I_{i,\ell} = 1$  if the  $\ell$ -th element  $[\text{vec}(\mathbf{W})]_\ell$  belongs to group  $i$  (i.e.,  $[\text{vec}(\mathbf{W})]_\ell$  is the element of  $\mathbf{w}_i$ ), otherwise  $I_{i,\ell} = 0$ .

## Propositions and lemmas used for deriving Theorem 1 in main paper

**Proposition 3** Under  $r_{ij} \geq 0$ ,  $r_{ij} = r_{ji}$ ,  $r_{ii} = 0$ , we have

$$\text{vec}(\mathbf{W})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}) = \sum_{i,j=1}^n r_{ij} \frac{\|\mathbf{w}_i - \mathbf{w}_j\|_2^2}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2},$$

where

$$\begin{aligned} \mathbf{F}_g^{(t)} &= \mathbf{I}_d \otimes \mathbf{C}^{(t)}, \\ [\mathbf{C}^{(t)}]_{i,j} &= \begin{cases} \sum_{j'=1}^n \frac{r_{ij'}}{\|\mathbf{w}_i^{(t)} - \mathbf{w}_{j'}^{(t)}\|_2} - \frac{r_{ij}}{\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} & (i = j) \\ \frac{-r_{ij}}{\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} & (i \neq j) \end{cases}. \end{aligned}$$

*Proof:*

$$\begin{aligned} &\sum_{i,j=1}^n r_{ij} \frac{\|\mathbf{w}_i - \mathbf{w}_j\|_2^2}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} \\ &= \sum_{i=1}^n \mathbf{w}_i^\top \mathbf{w}_i \sum_{j=1}^n \frac{r_{ij}}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} + \sum_{j=1}^n \mathbf{w}_j^\top \mathbf{w}_j \sum_{i=1}^n \frac{r_{ij}}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} - 2 \sum_{i=1}^n \sum_{j=1}^n \mathbf{w}_i^\top \mathbf{w}_j \frac{r_{ij}}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} \\ &= \text{tr}(\mathbf{W}^\top \mathbf{C}^{(t)} \mathbf{W}) \\ &= \text{vec}(\mathbf{W})^\top (\mathbf{I}_d \otimes \mathbf{C}^{(t)}) \text{vec}(\mathbf{W}), \end{aligned}$$

where  $\mathbf{I}_d \in \mathbb{R}^{d \times d}$  is the identity matrix,  $\text{tr}(\cdot)$  is the trace operator, and  $\text{vec}(\cdot)$  is the vectorization operator.

**Lemma 4** Under the updating rule of Eq. (6),

$$\tilde{J}(\mathbf{W}^{(t+1)}) - \tilde{J}(\mathbf{W}^{(t)}) \leq 0.$$

*Proof:* Under the updating rule of Eq. (6), since Eq.(5) is a convex function and the optimal solution is obtained by solving  $\frac{\partial \tilde{J}(\mathbf{W})}{\partial \mathbf{W}} = 0$ , the obtained solution  $\mathbf{W}^{(t+1)}$  is the global solution. That is,  $\tilde{J}(\mathbf{W}^{(t+1)}) \leq \tilde{J}(\mathbf{W}^{(t)})$ .

**Lemma 5** For any nonzero vectors  $\mathbf{w}, \mathbf{w}^{(t)} \in \mathbb{R}^d$ , the following inequality holds Nie et al. [2010]:

$$\|\mathbf{w}\|_2 - \frac{\|\mathbf{w}\|_2^2}{2\|\mathbf{w}^{(t)}\|_2} \leq \|\mathbf{w}^{(t)}\|_2 - \frac{\|\mathbf{w}^{(t)}\|_2^2}{2\|\mathbf{w}^{(t)}\|_2}.$$

**Lemma 6** For  $r_{i,j} \geq 0, \forall i, j$ , the following inequality holds for any non-zero vectors  $\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}, \mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}$ :

$$\begin{aligned} & \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}\|_2 - \text{vec}(\mathbf{W}^{(t+1)})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}^{(t+1)}) \\ & - \left( \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2 - \text{vec}(\mathbf{W}^{(t)})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}^{(t)}) \right) \leq 0. \end{aligned}$$

*Proof:*  $\text{vec}(\mathbf{W})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W})$  can be written as

$$\text{vec}(\mathbf{W})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}) = \sum_{i,j=1}^n r_{ij} \frac{\|\mathbf{w}_i - \mathbf{w}_j\|_2^2}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2}.$$

where  $r_{ij} \geq 0$ .

Then, the left hand side equation can be written as

$$\begin{aligned} \Delta_g &= \sum_{i,j=1}^n r_{ij} \left( \|\mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}\|_2 - \frac{\|\mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}\|_2^2}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} \right) \\ &\quad - \sum_{i,j=1}^n r_{ij} \left( \|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2 - \frac{\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2^2}{2\|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2} \right). \end{aligned}$$

Using Lemma 5,  $\Delta_g \leq 0$ . □

**Lemma 7** The following inequality holds for any non-zero vectors Kong et al. [2014]:

$$\begin{aligned} & \sum_{i=1}^n \|\mathbf{w}_i^{(t+1)}\|_1^2 - \text{vec}(\mathbf{W}^{(t+1)})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}^{(t+1)}) \\ & - \left( \sum_{i=1}^n \|\mathbf{w}_i^{(t)}\|_1^2 - \text{vec}(\mathbf{W}^{(t)})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}^{(t)}) \right) \leq 0. \end{aligned} \tag{1}$$

*Proof:*  $\text{vec}(\mathbf{W})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W})$  can be written as

$$\begin{aligned} \text{vec}(\mathbf{W})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}) &= \sum_{\ell=1}^{dn} [\text{vec}(\mathbf{W})]_\ell^2 \sum_{i=1}^n \frac{I_{i,\ell} \|\mathbf{w}_i^{(t)}\|_1}{[\text{vec}(|\mathbf{W}^{(t)}|)]_\ell} \\ &= \sum_{i=1}^n \left( \sum_{j=1}^d \frac{[\mathbf{w}_i^{(t)}]_j^2}{[|\mathbf{w}_i^{(t)}|]_j} \right) \|\mathbf{w}_i^{(t)}\|_1. \end{aligned}$$

Thus, the left hand equation is written as

$$\begin{aligned} \Delta_e &= \sum_{i=1}^n \left[ \left( \sum_{j=1}^d [\mathbf{w}_i^{(t+1)}]_j \right)^2 - \left( \sum_{j=1}^d \frac{[\mathbf{w}_i^{(t+1)}]_j^2}{[|\mathbf{w}_i^{(t)}|]_j} \right) \left( \sum_{j=1}^d [\mathbf{w}_i^{(t)}]_j \right) \right] \\ &= \sum_{i=1}^n \left[ \left( \sum_{j=1}^d a_j^{(t)} b_j^{(t)} \right)^2 - \left( \sum_{j=1}^d (a_j^{(t)})^2 \right) \left( \sum_{j=1}^d (b_j^{(t)})^2 \right) \right] \leq 0, \end{aligned}$$

where  $a_j^{(t)} = \frac{[\|\mathbf{w}_i^{(t+1)}\|]_j}{\sqrt{[\|\mathbf{w}_i^{(t)}\|]_j}}$  and  $b_j^{(t)} = \sqrt{[\|\mathbf{w}_i^{(t)}\|]_j}$ , and  $\text{vec}(\mathbf{W}^{(t)})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}^{(t)}) = \sum_{i=1}^n \|\mathbf{w}_i^{(t)}\|_1^2$ . The inequality holds due to cauchy inequality Steele [2004].  $\square$

**Lemma 8** For  $r_{i,j} \geq 0, \forall i, j$ , the following inequality holds for any non-zero vectors  $\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}, \mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}$ :

$$J(\mathbf{W}^{(t+1)}) - J(\mathbf{W}^{(t)}) \leq \tilde{J}(\mathbf{W}^{(t+1)}) - \tilde{J}(\mathbf{W}^{(t)}).$$

*Proof:* The difference between the right and left side equations is given as

$$\begin{aligned} \Delta &= J(\mathbf{W}^{(t+1)}) - J(\mathbf{W}^{(t)}) - (\tilde{J}(\mathbf{W}^{(t+1)}) - \tilde{J}(\mathbf{W}^{(t)})) \\ &= \lambda_1 \left( \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i^{(t+1)} - \mathbf{w}_j^{(t+1)}\|_2 - \text{vec}(\mathbf{W}^{(t+1)})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}^{(t+1)}) \right. \\ &\quad \left. - \left[ \sum_{i,j=1}^n r_{ij} \|\mathbf{w}_i^{(t)} - \mathbf{w}_j^{(t)}\|_2 - \text{vec}(\mathbf{W}^{(t)})^\top \mathbf{F}_g^{(t)} \text{vec}(\mathbf{W}^{(t)}) \right] \right) \\ &\quad + \lambda_2 \left( \sum_{i=1}^n \|\mathbf{w}_i^{(t+1)}\|_1^2 - \text{vec}(\mathbf{W}^{(t+1)})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}^{(t+1)}) \right. \\ &\quad \left. - \left[ \sum_{i=1}^n \|\mathbf{w}_i^{(t)}\|_1^2 - \text{vec}(\mathbf{W}^{(t)})^\top \mathbf{F}_e^{(t)} \text{vec}(\mathbf{W}^{(t)}) \right] \right) \end{aligned}$$

Based on Lemma 6 and 7,  $\Delta \leq 0$ .  $\square$

## References

- Deguang Kong, Ryohei Fujimaki, Ji Liu, Feiping Nie, and Chris Ding. Exclusive feature learning on arbitrary structures via  $\ell_{12}$ -norm. In *NIPS*, 2014.
- Feiping Nie, Heng Huang, Xiao Cai, and Chris H Ding. Efficient and robust feature selection via joint  $\ell_{2,1}$ -norms minimization. In *NIPS*, 2010.
- J Michael Steele. An introduction to the art of mathematical inequalities: The Cauchy-Schwarz master class, 2004.