

Title

Chao Cheng No.2017310132

February 3, 2021

Contents

1	Algorithm Details	2
1.1	Update of β	2
1.2	Update of μ	2
1.3	Update of z	3
1.4	Update of s	4
1.5	Update of w	6
1.6	Choice of $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$	8
1.7	Algorithm with L_2 Loss	11
A	Re-deduction of the Algorithm with More Parameters	13
A.1	Update steps for β and μ	13
A.1.1	Update β , Coordinate Descent	13
A.1.2	Update μ , Coordinate Descent	13
A.1.3	Update β and μ simultaneously	14
A.2	Update of z	14
A.3	Update of s	14
A.4	Update of w	15
A.5	Choice of $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$	16

1 Algorithm Details

In this section, we will provide the details of each upstating step of the algorithm.

1.1 Update of β

$$\begin{aligned}\beta^{(k+1)} &= \arg \min_{\beta} L(\beta, \mu^{(k)}, z^{(k)}, s^{(k)}, w^{(k)}, q_1^{(k)}, q_2^{(k)}, q_3^{(k)}) \\ &= \arg \min_{\beta} \left\{ \frac{r_1}{2} \|\mathbf{y} - \mu^{(k)} - \mathbf{X}\beta - \mathbf{z}^{(k)}\|_2^2 + \frac{r_3}{2} \|\beta - w^{(k)}\|_2^2 \right. \\ &\quad \left. + \langle \mathbf{y} - \mu^{(k)} - \mathbf{X}\beta - \mathbf{z}^{(k)}, \mathbf{q}_1^{(k)} \rangle + \langle \beta - w^{(k)}, \mathbf{q}_3^{(k)} \rangle \right\}\end{aligned}$$

This is quadratic in β and we can take the first derivative and set it to 0.

$$\begin{aligned}\mathbf{0}_p &= \frac{\partial}{\partial \beta} \left(\frac{r_1}{2} \|\mathbf{y} - \mu^{(k)} - \mathbf{X}\beta - \mathbf{z}^{(k)}\|_2^2 + \frac{r_3}{2} \|\beta - w^{(k)}\|_2^2 - \beta^T \mathbf{X}^T \mathbf{q}_1^{(k)} + \beta^T \mathbf{q}_3^{(k)} \right) \\ &= \frac{\partial}{\partial \beta} \left(\beta^T \left(\frac{r_1}{2} \mathbf{X}^T \mathbf{X} + \frac{r_3}{2} \mathbf{I}_p \right) \beta + \beta^T (-r_1 \mathbf{X}^T (\mathbf{y} - \mu - \mathbf{z}) - r_3 w - \mathbf{X}^T \mathbf{q}_1 + \mathbf{q}_3) \right) \\ &= (r_1 \mathbf{X}^T \mathbf{X} + r_3 \mathbf{I}_p) \beta + (-r_1 \mathbf{X}^T (\mathbf{y} - \mu - \mathbf{z}) - r_3 w - \mathbf{X}^T \mathbf{q}_1 + \mathbf{q}_3)\end{aligned}$$

Then the solution is

$$\beta^{(k+1)} = (r_1 \mathbf{X}^T \mathbf{X} + r_3 \mathbf{I}_p)^{-1} (r_1 \mathbf{X}^T (\mathbf{y} - \mu - \mathbf{z}) + r_3 w + \mathbf{X}^T \mathbf{q}_1 - \mathbf{q}_3)$$

In this solution we have to solve the inverse of a $p \times p$ matrix. But if $p > n$ or even $p \gg n$, this would require a lot of computational resources. Nevertheless, if we left multiply \mathbf{X} on both side of the first derivative equation, then by some algebra we can show that

$$\beta^{(k+1)} = \frac{1}{r_3} \left[\mathbf{I}_p - r_1 \mathbf{X}^T (r_1 \mathbf{X} \mathbf{X}^T + r_3 \mathbf{I}_n)^{-1} \mathbf{X} \right] (r_1 \mathbf{X}^T (\mathbf{y} - \mu - \mathbf{z}) + r_3 w + \mathbf{X}^T \mathbf{q}_1 - \mathbf{q}_3)$$

which only requires to solve the inverse of a $n \times n$ matrix.

1.2 Update of μ

Like in 1.1, we can write

$$\begin{aligned}\mu^{(k+1)} &= \arg \min_{\mu} L(\beta, \mu, z, s, w, q_1, q_2, q_3) \\ &= \arg \min_{\mu} \left\{ \frac{r_1}{2} \|\mathbf{y} - \mu - \mathbf{X}\beta - \mathbf{z}\|_2^2 + \frac{r_2}{2} \|\mathbf{D}\mu - \mathbf{s}\|_2^2 \right. \\ &\quad \left. + \langle \mathbf{y} - \mu - \mathbf{X}\beta - \mathbf{z}, \mathbf{q}_1 \rangle + \langle \mathbf{D}\mu - \mathbf{s}, \mathbf{q}_2 \rangle \right\}\end{aligned}$$

and take the first derivative then set it to 0

$$\begin{aligned}0 &= \frac{\partial}{\partial \mu} \left(\frac{r_1}{2} \|\mathbf{y} - \mu - \mathbf{X}\beta - \mathbf{z}\|_2^2 + \frac{r_2}{2} \|\mathbf{D}\mu - \mathbf{s}\|_2^2 - \mu^T \mathbf{q}_1 + \mu^T \mathbf{D}^T \mathbf{q}_2 \right) \\ &= \frac{\partial}{\partial \mu} \left(\mu^T \left(\frac{r_1}{2} \mathbf{I}_n + \frac{r_2}{2} \mathbf{D}^T \mathbf{D} \right) \mu + \mu^T (-r_1 (\mathbf{y} - \mathbf{X}\beta - \mathbf{z}) - r_2 \mathbf{D}^T \mathbf{s} - \mathbf{q}_1 + \mathbf{D}^T \mathbf{q}_2) \right) \\ &= (r_1 \mathbf{I}_n + r_2 \mathbf{D}^T \mathbf{D}) \mu + (-r_1 (\mathbf{y} - \mathbf{X}\beta - \mathbf{z}) - r_2 \mathbf{D}^T \mathbf{s} - \mathbf{q}_1 + \mathbf{D}^T \mathbf{q}_2)\end{aligned}$$

So the solution is

$$\mu = (r_1 \mathbf{I}_n + r_2 \mathbf{D}^T \mathbf{D})^{-1} (r_1 (\mathbf{y} - \mathbf{X}\beta - \mathbf{z}) + r_2 \mathbf{D}^T \mathbf{s} + \mathbf{q}_1 - \mathbf{D}^T \mathbf{q}_2)$$

1.3 Update of \mathbf{z}

Here, we use L-1 loss as an example.

$$\begin{aligned} \mathbf{z}^{(k+1)} &= \arg \min_{\mathbf{z}} L(\boldsymbol{\beta}, \boldsymbol{\mu}, \mathbf{z}, \mathbf{s}, \mathbf{w}, \mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3) \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{n} \sum_{i=1}^n |z_i| + \frac{r_1}{2} \|\mathbf{y} - \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}\|_2^2 + \langle \mathbf{y} - \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}, \mathbf{q}_1 \rangle \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{n} \sum_{i=1}^n |z_i| + \frac{r_1}{2} (\mathbf{z}^T \mathbf{z} - 2\mathbf{z}^T (\mathbf{y} - \boldsymbol{\mu} - \mathbf{X}^T \boldsymbol{\beta})) - \mathbf{z}^T \mathbf{q}_1 \right\} \end{aligned}$$

Note that this problem can be solved coordinate-wisely by

$$z_i^{(k+1)} = \arg \min_{z_i} \frac{1}{n} |z_i| + \frac{r_1}{2} (z_i^2 - 2z_i (y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta})) - z_i q_{1,i}$$

By the subgradient method, we have

$$\begin{aligned} 0 &\in \partial \left(\frac{1}{n} |z_i| + \frac{r_1}{2} (z_i^2 - 2z_i (y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta})) - z_i q_{1,i} \right) \Big|_{z_i = z_i^{(k+1)}} \\ \implies 0 &\in \frac{\partial (|z_i|)}{n} \Big|_{z_i = z_i^{(k+1)}} + r_1 \left(z_i^{(k+1)} - y_i + \mu_i + \mathbf{x}_i^T \boldsymbol{\beta} - \frac{q_{1,i}}{r_1} \right) \end{aligned}$$

where

$$\partial (|z_i|) \Big|_{z_i = z_i^{(k+1)}} = \begin{cases} 1 & z_i^{(k+1)} > 0 \\ -1 & z_i^{(k+1)} < 0 \\ [-1, 1] & z_i^{(k+1)} = 0 \end{cases}$$

Therefore

$$\begin{cases} 0 = \frac{1}{n} + r_1 \left(z_i^{(k+1)} - y_i + \mu_i + \mathbf{x}_i^T \boldsymbol{\beta} - \frac{q_{1,i}}{r_1} \right) & z_i^{(k+1)} > 0 \\ 0 = -\frac{1}{n} + r_1 \left(z_i^{(k+1)} - y_i + \mu_i + \mathbf{x}_i^T \boldsymbol{\beta} - \frac{q_{1,i}}{r_1} \right) & z_i^{(k+1)} < 0 \\ 0 \in \left[-\frac{1}{n}, \frac{1}{n} \right] + r_1 \left(z_i^{(k+1)} - y_i + \mu_i + \mathbf{x}_i^T \boldsymbol{\beta} - \frac{q_{1,i}}{r_1} \right) & z_i^{(k+1)} = 0 \end{cases}$$

Hence

$$z_i^{(k+1)} = \begin{cases} y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} - \frac{1}{nr_1} & \frac{1}{nr_1} < y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} \\ 0 & -\frac{1}{nr_1} \leq y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} \leq \frac{1}{nr_1} \\ y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} + \frac{1}{nr_1} & y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} < -\frac{1}{nr_1} \end{cases}$$

Here we introduce the soft-thresholding function $S(x, \lambda)$ for $x \in \mathcal{R}$ and $\lambda > 0$:

$$S(x, \lambda) = \begin{cases} x - \lambda & \lambda < x \\ 0 & |x| \leq \lambda \\ x + \lambda & x < -\lambda \end{cases}$$

Then we can summary the update of $z_i^{(k+1)}$ by

$$z_i^{(k+1)} = S \left(y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1}, \frac{1}{nr_1} \right)$$

Note: If we use the **Huber loss with parameter c** instead of the absolute value loss, then similarly the update of $z_i^{(k+1)}$ would be

$$z_i^{(k+1)} = \begin{cases} y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} - \frac{c}{nr_1} & \frac{c}{nr_1} + c < y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} \\ \frac{y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1}}{1 + \frac{1}{nr_1}} & -\frac{c}{nr_1} - c \leq y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} \leq \frac{c}{nr_1} + c \\ y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} + \frac{c}{nr_1} & y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{q_{1,i}}{r_1} < -\frac{c}{nr_1} - c \end{cases}$$

FYI: The huber loss with parameter c is defined as

$$\rho(x; c) = \begin{cases} \frac{1}{2}x^2 & |x| \leq c \\ c|x| - \frac{1}{2}c^2 & |x| > c \end{cases}$$

1.4 Update of \mathbf{s}

Here we only consider the SCAD penalty case, which means $P_{\lambda_1}(x) = P_{\lambda_1, \gamma_1}(x)$ where γ_1 is a SCAD parameter.

$$\begin{aligned} \mathbf{s}^{(k+1)} &= \arg \min_{\mathbf{s}} L(\boldsymbol{\beta}, \boldsymbol{\mu}, \mathbf{z}, \mathbf{s}, \mathbf{w}, \mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3) \\ &= \arg \min_{\mathbf{s}} \left\{ \sum_{1 \leq i < j \leq n} P_{\lambda_1, \gamma_1}(s_{ij}) + \frac{r_2}{2} \|\mathbf{D}\boldsymbol{\mu} - \mathbf{s}\|_2^2 + \langle \mathbf{D}\boldsymbol{\mu} - \mathbf{s} \rangle \right\} \\ &= \arg \min_{\mathbf{s}} \left\{ \sum_{1 \leq i < j \leq n} P_{\lambda_1, \gamma_1}(s_{ij}) + \frac{r_2}{2} (\mathbf{s}^T \mathbf{s} - 2\mathbf{s}^T \mathbf{D}\boldsymbol{\mu}) - \mathbf{s}^T \mathbf{q}_2 \right\} \end{aligned}$$

Note that $\{s_{ij}\}$ are actually mutually independent in this optimization problem. So like in 1.3, this can also be solved coordinate-wisely. And we try to solve

$$s_{ij}^{(k+1)} = \arg \min_{s_{ij}} \left\{ P_{\lambda_1, \gamma_1}(s_{ij}) + \frac{r_2}{2} (s_{ij}^2 - 2s_{ij}(\mu_i - \mu_j)) - q_{2,ij}s_{ij} \right\}$$

and likewise

$$\begin{aligned} 0 &\in \partial \left(P_{\lambda_1, \gamma_1}(s_{ij}) + \frac{r_2}{2} (s_{ij}^2 - 2s_{ij}(\mu_i - \mu_j)) - q_{2,ij}s_{ij} \right) \Big|_{s_{ij}=s_{ij}^{(k+1)}} \\ \implies 0 &\in \partial P_{\lambda_1, \gamma_1}(s_{ij}) \Big|_{s_{ij}=s_{ij}^{(k+1)}} + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) \end{aligned}$$

For the SCAD penalty, we have

$$P'_{\lambda_1, \gamma_1}(s_{ij}) = \lambda_1 \left\{ I(s_{ij} \leq \lambda_1) + \frac{(\gamma_1 \lambda_1 - s_{ij})_+}{(\gamma_1 - 1) \lambda_1} I(s_{ij} > \lambda_1) \right\}$$

for some $\gamma_1 > 2$ and $s_{ij} > 0$. And $P'_{\lambda_1, \gamma_1}(s_{ij}) = -P'_{\lambda_1, \gamma_1}(-s_{ij})$ when $s_{ij} < 0$. Therefore

$$\left\{ \begin{array}{ll} 0 = 0 + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & \gamma_1 \lambda_1 < s_{ij} \\ 0 = \frac{\gamma_1 \lambda_1 - s_{ij}}{\gamma_1 - 1} + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & \lambda_1 < s_{ij} \leq \gamma_1 \lambda_1 \\ 0 = \lambda_1 + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & 0 < s_{ij} \leq \lambda_1 \\ 0 \in [-\lambda_1, \lambda_1] + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & s_{ij} = 0 \\ 0 = -\lambda_1 + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & -\lambda_1 \leq s_{ij} < 0 \\ 0 = -\frac{\gamma_1 \lambda_1 + s_{ij}}{\gamma_1 - 1} + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & -\gamma_1 \lambda_1 \leq s_{ij} < -\lambda_1 \\ 0 = 0 + r_2 \left(s_{ij}^{(k+1)} - (\mu_i - \mu_j) - \frac{q_{2,ij}}{r_2} \right) & s_{ij} \leq -\gamma_1 \lambda_1 \end{array} \right.$$

Hence we have

$$s_{ij}^{(k+1)} = \left\{ \begin{array}{ll} \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \gamma_1 \lambda_1 < \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \\ \frac{r_2(\gamma_1 - 1) \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right) - \gamma_1 \lambda_1}{r_2(\gamma_1 - 1) - 1} & \left(1 + \frac{1}{r_2} \right) \lambda_1 < \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \leq \gamma_1 \lambda_1 \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} - \frac{\lambda_1}{r_2} & \frac{\lambda_1}{r_2} < \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \leq \left(1 + \frac{1}{r_2} \right) \lambda_1 \\ 0 & -\frac{\lambda_1}{r_2} \leq \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \leq \frac{\lambda_1}{r_2} \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} + \frac{\lambda_1}{r_2} & -\left(1 + \frac{1}{r_2} \right) \lambda_1 \leq \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} < -\frac{\lambda_1}{r_2} \\ \frac{r_2(\gamma_1 - 1) \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right) + \gamma_1 \lambda_1}{r_2(\gamma_1 - 1) - 1} & -\gamma_1 \lambda_1 \leq \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} < -\left(1 + \frac{1}{r_2} \right) \lambda_1 \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \leq -\gamma_1 \lambda_1 \end{array} \right.$$

Note: During the deduction, condition $r_2(\gamma_1 - 1) > 1$ (i.e. $r_2 > \frac{1}{\gamma_1 - 1}$) is needed to guarantee the result. **Also with this condition, the objective function is convex in s_{ij} .**

This update can be summarized with the help of the soft-thresholding function

$$s_{ij}^{(k+1)} = \left\{ \begin{array}{ll} S \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{\lambda_1}{r_2} \right) & \left| \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right| \leq \left(1 + \frac{1}{r_2} \right) \lambda_1 \\ \frac{S \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{\gamma_1 \lambda_1}{r_2(\gamma_1 - 1)} \right)}{1 - \frac{1}{r_2(\gamma_1 - 1)}} & \left(1 + \frac{1}{r_2} \right) \lambda_1 < \left| \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right| \leq \gamma_1 \lambda_1 \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \left| \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right| > \gamma_1 \lambda_1 \end{array} \right.$$

For MCP or LASSO, the technique are quite similar. (And the result would be shown below:) If the penalty is Lasso with parameter λ_1 then

$$s_{ij}^{(k+1)} = S \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{\lambda_1}{r_2} \right)$$

If the penalty is MCP with parameter λ_1 and γ_1 , then the panalty takes the form

$$P_{\lambda_1, \gamma_1}(s_{ij}) = \lambda_1 \int_0^{s_{ij}} \left(1 - \frac{x}{\lambda_1 \gamma_1} \right)_+ dx$$

for some $\gamma_1 > 1$ and $s_{ij} \geq 0$. And $P'_{\lambda_1, \gamma_1}(s_{ij}) = -P'_{\lambda_1, \gamma_1}(-s_{ij})$ when $s_{ij} < 0$. Therefore the solution is

$$s_{ij}^{(k+1)} = \begin{cases} \frac{S \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{\lambda_1}{r_2} \right)}{1 - \frac{1}{r_2 \gamma_1}} & \left| \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right| \leq \gamma_1 \lambda_1 \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \left| \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} \right| > \gamma_1 \lambda_1 \end{cases}$$

Note: Here the condition $r_2 \gamma_1 > 1$ (i.e. $r_2 > \frac{1}{\gamma_1}$) is needed, and the objective function is convex in s_{ij} .

1.5 Update of \mathbf{w}

Like in 1.4, here we consider SCAD penalty as an example. Therefore $P_{\lambda_2}(x) = P_{\lambda_2, \gamma_2}(x)$ where γ_2 is a SCAD parameter.

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \arg \min_{\mathbf{w}} L(\boldsymbol{\beta}, \boldsymbol{\mu}, \mathbf{z}, \mathbf{s}, \mathbf{w}, \mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3) \\ &= \arg \min_{\mathbf{w}} \left\{ \sum_{j=1}^p P_{\lambda_2, \gamma_2}(w_j) + \frac{r_3}{2} \|\boldsymbol{\beta} - \mathbf{w}\|_2^2 + \langle \boldsymbol{\beta} - \mathbf{w}, \mathbf{q}_3 \rangle \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ \sum_{j=1}^p P_{\lambda_2, \gamma_2}(w_j) + \frac{r_3}{2} (\mathbf{w}^T \mathbf{w} - 2\mathbf{w}^T \boldsymbol{\beta}) - \mathbf{w}^T \mathbf{q}_3 \right\} \end{aligned}$$

Again, this can be solved elementwisely by

$$w_j^{(k+1)} = \arg \min_{w_j} \left\{ P_{\lambda_2, \gamma_2}(w_j) + \frac{r_3}{2} (w_j^2 - 2w_j \beta_j) - w_j q_{3,j} \right\}$$

Like before, subgradient method provides us

$$\begin{aligned} 0 &\in \partial \left(P_{\lambda_2, \gamma_2}(w_j) + \frac{r_3}{2} (w_j^2 - 2w_j \beta_j) - w_j q_{3,j} \right) \Big|_{w_j = w_j^{(k+1)}} \\ \implies 0 &\in \partial P_{\lambda_2, \gamma_2}(w_j) \Big|_{w_j = w_j^{(k+1)}} + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) \end{aligned}$$

which means

$$\left\{ \begin{array}{ll} 0 = 0 + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & \gamma_2 \lambda_2 < w_j \\ 0 = \frac{\gamma_2 \lambda_2 - w_j}{\gamma_2 - 1} + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & \lambda_2 < w_j \leq \gamma_2 \lambda_2 \\ 0 = \lambda_2 + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & 0 < w_j \leq \lambda_2 \\ 0 \in [-\lambda_2, \lambda_2] + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & w_j = 0 \\ 0 = -\lambda_2 + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & -\lambda_2 \leq w_j < 0 \\ 0 = -\frac{\lambda_2 \gamma_2 + w_j}{\gamma_2 - 1} + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & -\gamma_2 \lambda_2 \leq w_j < -\lambda_2 \\ 0 = 0 + r_3 \left(w_j - \beta_j - \frac{q_{3,j}}{r_3} \right) & w_j \leq -\gamma_2 \lambda_2 \end{array} \right.$$

Therefore

$$w_j^{(k+1)} = \left\{ \begin{array}{ll} \beta_j + \frac{q_{3,j}}{r_3} & \gamma_2 \lambda_2 < \beta_j + \frac{q_{3,j}}{r_3} \\ \frac{r_3 (\gamma_2 - 1) \left(\beta_j + \frac{q_{3,j}}{r_3} \right) - \gamma_2 \lambda_2}{r_3 (\gamma_2 - 1) - 1} & \left(1 + \frac{1}{r_3} \right) \lambda_2 < \beta_j + \frac{q_{3,j}}{r_3} \leq \gamma_2 \lambda_2 \\ \beta_j + \frac{q_{3,j}}{r_3} - \frac{\lambda_2}{r_3} & \frac{\lambda_2}{r_3} < \beta_j + \frac{q_{3,j}}{r_3} \leq \left(1 + \frac{1}{r_3} \right) \lambda_2 \\ 0 & -\frac{\lambda_2}{r_3} \leq \beta_j + \frac{q_{3,j}}{r_3} \leq \frac{\lambda_2}{r_3} \\ \beta_j + \frac{q_{3,j}}{r_3} + \frac{\lambda_2}{r_3} & -\left(1 + \frac{1}{r_3} \right) \lambda_2 \leq \beta_j + \frac{q_{3,j}}{r_3} < -\frac{\lambda_2}{r_3} \\ \frac{r_3 (\gamma_2 - 1) \left(\beta_j + \frac{q_{3,j}}{r_3} \right) + \lambda_2 \gamma_2}{r_3 (\gamma_2 - 1) - 1} & -\gamma_2 \lambda_2 \leq \beta_j + \frac{q_{3,j}}{r_3} < -\left(1 + \frac{1}{r_3} \right) \lambda_2 \\ \beta_j + \frac{q_{3,j}}{r_3} & \beta_j + \frac{q_{3,j}}{r_3} < -\gamma_2 \lambda_2 \end{array} \right.$$

Note: Here we need $r_3 (\gamma_2 - 1) > 1$ (i.e. $r_3 > \frac{1}{\gamma_2 - 1}$) to guarantee the result.

We can summary the result of $w_j^{(k+1)}$ as

$$w_j^{(k+1)} = \left\{ \begin{array}{ll} S \left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{\lambda_2}{r_3} \right) & \left| \beta_j + \frac{q_{3,j}}{r_3} \right| \leq \left(1 + \frac{1}{r_3} \right) \lambda_2 \\ \frac{S \left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{\gamma_2 \lambda_2}{r_3 (\gamma_2 - 1)} \right)}{1 - \frac{1}{r_3 (\gamma_2 - 1)}} & \left(1 + \frac{1}{r_3} \right) \lambda_2 < \left| \beta_j + \frac{q_{3,j}}{r_3} \right| \leq \gamma_2 \lambda_2 \\ \beta_j + \frac{q_{3,j}}{r_3} & \left| \beta_j + \frac{q_{3,j}}{r_3} \right| > \gamma_2 \lambda_2 \end{array} \right.$$

Results for MCP and Lasso are presented as following. If the penalty is Lasso with parameter λ_2 then

$$w_j^{(k+1)} = S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{\lambda_2}{r_3}\right)$$

If the penalty is MCP with parameter λ_2 and γ_2 , then

$$w_j^{(k+1)} = \begin{cases} \frac{S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{\lambda_2}{r_3}\right)}{1 - \frac{1}{r_3\gamma_2}} & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| \leq \gamma_2\lambda_2 \\ \beta_j + \frac{q_{3,j}}{r_3} & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| > \gamma_2\lambda_2 \end{cases}$$

1.6 Choice of $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$

Need to be double checked.

Apparently for some large enough $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$, β and μ have to be shrinked to $\mathbf{0}_p$ and \mathbf{c}_n respectively, where $\mathbf{0}_p$ is a length- p zero vector and \mathbf{c}_n is a length- n constant vector. Actually, it's clear that all elements of \mathbf{c}_n have to be the medium of $\{y_i, i = 1, 2, \dots, n\}$. Therefore we can write

$$\begin{pmatrix} \mathbf{c}_n \\ \mathbf{0}_p \end{pmatrix} \in \arg \min_{\mu, \beta} \frac{1}{n} \sum_{i=1}^n |y_i - \mu_i - \mathbf{x}_i^T \beta| + \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) + \sum_{j=1}^p P_{\lambda_2}(\beta_j)$$

which means

$$\begin{pmatrix} \mathbf{0}_n \\ \mathbf{0}_p \end{pmatrix} \in \partial \left\{ \frac{1}{n} \sum_{i=1}^n |y_i - \mu_i - \mathbf{x}_i^T \beta| + \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) + \sum_{j=1}^p P_{\lambda_2}(\beta_j) \right\} \bigg|_{\begin{pmatrix} \mu \\ \beta \end{pmatrix} = \begin{pmatrix} \mathbf{c}_n \\ \mathbf{0}_p \end{pmatrix}}$$

That is to say

$$\begin{pmatrix} \mathbf{0}_n \\ \mathbf{0}_p \end{pmatrix} = \begin{pmatrix} \mathbf{d}_{1,n} \\ \mathbf{d}_{1,p} \end{pmatrix} + \begin{pmatrix} \mathbf{d}_{2,n} \\ \mathbf{0}_p \end{pmatrix} + \begin{pmatrix} \mathbf{0}_n \\ \mathbf{d}_{3,p} \end{pmatrix}$$

where

$$\begin{pmatrix} \mathbf{d}_{1,n} \\ \mathbf{d}_{1,p} \end{pmatrix} \in \partial \left\{ \frac{1}{n} \sum_{i=1}^n |y_{it} - \mu_i - \mathbf{x}_{it}^T \beta| \right\} \bigg|_{\begin{pmatrix} \mu \\ \beta \end{pmatrix} = \begin{pmatrix} \mathbf{c}_n \\ \mathbf{0}_p \end{pmatrix}} = \frac{1}{n} \begin{pmatrix} \partial |y_1 - c| \\ \vdots \\ \partial |y_n - c| \\ \sum_{i=1}^n (\partial |y_i - c|) \mathbf{x}_i \end{pmatrix}$$

$$\mathbf{d}_{2,n} \in \partial \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) \bigg|_{\mu = \mathbf{c}_n}$$

and

$$\mathbf{d}_{3,p} \in \partial \sum_{j=1}^p P_{\lambda_2}(\beta_j) \bigg|_{\beta = \mathbf{0}_p}$$

Note that

$$\begin{aligned}
\mathbf{0}_p &\in -\mathbf{d}_{3,p} + \partial \sum_{j=1}^p P_{\lambda_2}(\hat{\beta}_j) && (\text{where } \hat{\beta} = \mathbf{0}_p) \\
\Rightarrow \mathbf{0}_p &\in \hat{\beta} - (\hat{\beta} + \mathbf{d}_{3,p}) + \partial \sum_{j=1}^p P_{\lambda_2}(\hat{\beta}_j) \\
\Rightarrow \mathbf{0}_p &\in \partial \left\{ \frac{1}{2} \left\| \beta - (\hat{\beta} + \mathbf{d}_{3,p}) \right\|_2^2 + \sum_{j=1}^p P_{\lambda_2}(\beta_j) \right\} \Big|_{\beta=\hat{\beta}} \\
\Rightarrow \hat{\beta} &\in \arg \min_{\beta} \frac{1}{2} \left\| \beta - (\hat{\beta} + \mathbf{d}_{3,p}) \right\|_2^2 + \sum_{j=1}^p P_{\lambda_2}(\beta_j)
\end{aligned}$$

Again, we use SCAD penalty with parameter γ_2 as an example here, and we can solve this penalized OLS by

$$\mathbf{0}_p = \hat{\beta} = T_{\lambda_2, \gamma_2}^{SCAD}(\hat{\beta} + \mathbf{d}_{3,p}) = T_{\lambda_2, \gamma_2}^{SCAD}(\mathbf{d}_{3,p})$$

and by property of the SCAD thresholding function $T_{\lambda_2, \gamma_2}^{SCAD}$, we have

$$\begin{aligned}
\|\mathbf{d}_{3,p}\|_{\infty} &\leq \lambda_2 \\
\Rightarrow \|\mathbf{d}_{1,p}\|_{\infty} &\leq \lambda_2 \quad (\text{since } \mathbf{0}_p = \mathbf{d}_{1,p} + \mathbf{d}_{3,p})
\end{aligned}$$

Note that

$$\mathbf{d}_{1,p} = \frac{1}{n} \sum_{i=1}^n (\partial |y_i - c|) \mathbf{x}_{it}$$

Therefore

$$\lambda_2^{(0)} \geq \frac{\|\mathbf{d}_{1,p}\|_{\infty}}{n} = \frac{1}{n} \left\| \sum_{i=1}^n (\partial |y_i - c|) \mathbf{x}_i \right\|_{\infty}$$

Choose of $\lambda_1^{(0)}$:

$$\begin{aligned}
\mathbf{d}_{1,n} &\in \frac{1}{n} \begin{pmatrix} \partial |y_1 - c| \\ \partial |y_2 - c| \\ \vdots \\ \partial |y_n - c| \end{pmatrix} \\
\mathbf{d}_{2,n} &\in \partial \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) \Big|_{\mu=c_n}
\end{aligned}$$

and

$$\mathbf{d}_{1,n} + \mathbf{d}_{2,n} = \mathbf{0}_n$$

By the chain rule we have

$$\begin{aligned}
& \left. \partial \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) \right|_{\mu = \mathbf{c}_n} \\
&= \mathbf{D}^T \partial \sum_{1 \leq i < j \leq n} P_{\lambda_1}(s_{ij}) \Big|_{\mathbf{s} = \mathbf{0}} \\
&\in \mathbf{D}^T \begin{pmatrix} [-\lambda_1, \lambda_1] \\ [-\lambda_1, \lambda_1] \\ \vdots \\ [-\lambda_1, \lambda_1] \end{pmatrix}_{\frac{n(n-1)}{2} \times 1}
\end{aligned}$$

Therefore we can deduce

$$\begin{aligned}
\|\mathbf{d}_{2,n}\|_\infty &\in \left\| \mathbf{D}^T \begin{pmatrix} [-\lambda_1, \lambda_1] \\ [-\lambda_1, \lambda_1] \\ \vdots \\ [-\lambda_1, \lambda_1] \end{pmatrix} \right\|_\infty \leq (n-1) \lambda_1 \\
\Rightarrow \|\mathbf{d}_{1,n}\|_\infty &\leq (n-1) \lambda_1 \quad (\text{since } \mathbf{0}_n = \mathbf{d}_{1,n} + \mathbf{d}_{2,n})
\end{aligned}$$

Note that

$$\mathbf{d}_{1,n} \in \frac{1}{n} \begin{pmatrix} \partial |y_1 - c| \\ \partial |y_2 - c| \\ \vdots \\ \partial |y_n - c| \end{pmatrix}$$

Therefore

$$\lambda_1^{(0)} \geq \frac{\|\mathbf{d}_{1,n}\|_\infty}{n(n-1)} = \frac{1}{n(n-1)} \left\| \begin{pmatrix} \partial |y_1 - c| \\ \partial |y_2 - c| \\ \vdots \\ \partial |y_n - c| \end{pmatrix} \right\|_\infty$$

and a sufficient choice would be

$$\lambda_1^{(0)} = \frac{1}{n(n-1)}$$

Another way to look at $\lambda_1^{(0)}$:

Note that

$$\mathbf{d}_{1,n} + \mathbf{d}_{2,n} = \mathbf{0}_n,$$

$$\mathbf{d}_{2,n} \in \mathbf{D}^T \begin{pmatrix} [-\lambda_1, \lambda_1] \\ [-\lambda_1, \lambda_1] \\ \vdots \\ [-\lambda_1, \lambda_1] \end{pmatrix}_{\frac{n(n-1)}{2} \times 1} \quad \text{and} \quad \mathbf{d}_{1,n} \in \frac{-1}{n} \begin{pmatrix} \partial |y_1 - c| \\ \partial |y_2 - c| \\ \vdots \\ \partial |y_n - c| \end{pmatrix}.$$

Also by convex optimization theory, we have $\mathbf{1}^T \mathbf{d}_{1,n} = 0$. Therefore there exists a n -dimensional vector $\boldsymbol{\theta}$ such that

$$(\mathbf{D}^T \mathbf{D}) \boldsymbol{\theta} = -\mathbf{d}_{1,n},$$

since $\text{rank}(\mathbf{D}^T \mathbf{D}) = n - 1$. Hence one possible solution of $\mathbf{d}_{2,n}$ is

$$\mathbf{d}_{2,n} = -\mathbf{D}^T \mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{d}_{1,n},$$

as long as

$$\left\| \mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{d}_{1,n} \right\|_{\infty} \leq \lambda_1,$$

where $(\mathbf{D}^T \mathbf{D})^{-1}$ is the Moore-Penrose generalized inverse of $\mathbf{D}^T \mathbf{D}$. And the choice of $\lambda_1^{(0)}$ is

$$\lambda_1^{(0)} = \left\| \mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{d}_{1,n} \right\|_{\infty}$$

Note that it's easy to verify

$$(\mathbf{D}^T \mathbf{D})^{-1} = \begin{pmatrix} \frac{n-1}{n^2} & \frac{-1}{n^2} & \cdots & \frac{-1}{n^2} \\ \frac{-1}{n^2} & \frac{n-1}{n^2} & \cdots & \frac{-1}{n^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{-1}{n^2} & \cdots & \frac{-1}{n^2} & \frac{n-1}{n^2} \end{pmatrix},$$

and there are two and only two entries in each row of $\mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1}$ with their values being $\frac{1}{n}$ and $\frac{-1}{n}$. Hence we know that

$$\left\| \mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{d}_{1,n} \right\|_{\infty} \leq \frac{2}{n} \|\mathbf{d}_{1,n}\|_{\infty} = \frac{2}{n^2} \left\| \begin{pmatrix} \partial |y_1 - c| \\ \partial |y_2 - c| \\ \vdots \\ \partial |y_n - c| \end{pmatrix} \right\|_{\infty},$$

which is a less computation intensive bound for $\lambda_1^{(0)}$.

1.7 Algorithm with L_2 Loss

If we use L_2 loss, then there is no need for the \mathbf{z} part in the algorithm and the augmented lagrangian form would be

$$\begin{aligned} & L(\boldsymbol{\beta}, \boldsymbol{\mu}, \mathbf{s}, \mathbf{w}, \mathbf{q}_2, \mathbf{q}_3) \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \sum_{1 \leq i < j \leq n} P_{\lambda_1}(s_{ij}) + \sum_{j=1}^p P_{\lambda_2}(w_j) \\ &+ \frac{r_2}{2} \|\mathbf{D}\boldsymbol{\mu} - \mathbf{s}\|_2^2 + \frac{r_3}{2} \|\boldsymbol{\beta} - \mathbf{w}\|_2^2 + \langle \mathbf{D}\boldsymbol{\mu} - \mathbf{s}, \mathbf{q}_2 \rangle + \langle \boldsymbol{\beta} - \mathbf{w}, \mathbf{q}_3 \rangle \end{aligned}$$

The update of \mathbf{s} and \mathbf{w} is just the same as before. And the update of $\boldsymbol{\beta}$ and $\boldsymbol{\mu}$ is just to find the minimum of

$$\begin{aligned} & \boldsymbol{\beta}^T \left(\frac{1}{n} \mathbf{X}^T \mathbf{X} + \frac{r_3}{2} \mathbf{I}_p \right) \boldsymbol{\beta} + \boldsymbol{\mu}^T \left(\frac{1}{n} \mathbf{I}_n + \frac{r_2}{2} \mathbf{D}^T \mathbf{D} \right) \boldsymbol{\mu} + 2\boldsymbol{\beta}^T \mathbf{X}^T \boldsymbol{\mu} \\ & - 2\mathbf{Y}^T \mathbf{X} \boldsymbol{\beta} - r_3 \mathbf{w}^T \boldsymbol{\beta} + \mathbf{q}_3^T \boldsymbol{\beta} - 2\mathbf{Y}^T \boldsymbol{\mu} - r_2 \mathbf{s}^T \mathbf{D} \boldsymbol{\mu} + \mathbf{q}_2^T \mathbf{D}^T \boldsymbol{\mu} \end{aligned}$$

Taking subgradient and set it to $\mathbf{0}$ we can find the linear system is

$$\begin{pmatrix} \frac{1}{n} \mathbf{I}_n + \frac{r_2}{2} \mathbf{D}^T \mathbf{D} & \frac{1}{n} \mathbf{X} \\ \frac{1}{n} \mathbf{X}^T & \frac{1}{n} \mathbf{X}^T \mathbf{X} + \frac{r_3}{2} \mathbf{I}_p \end{pmatrix} \begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\beta} \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \mathbf{Y} + \frac{r_2}{2} \mathbf{D}^T \mathbf{s} - \frac{1}{2} \mathbf{D}^T \mathbf{q}_2 \\ \frac{1}{n} \mathbf{X}^T \mathbf{Y} + \frac{r_3}{2} \mathbf{w} - \frac{1}{2} \mathbf{q}_3 \end{pmatrix}$$

If the dimension is not large, we can directly find the update by

$$\begin{pmatrix} \boldsymbol{\mu}^{(k+1)} \\ \boldsymbol{\beta}^{(k+1)} \end{pmatrix} = \begin{pmatrix} \frac{1}{n}\mathbf{I}_n + \frac{r_2}{2}\mathbf{D}^T\mathbf{D} & \frac{1}{n}\mathbf{X} \\ \frac{1}{n}\mathbf{X}^T & \frac{1}{n}\mathbf{X}^T\mathbf{X} + \frac{r_3}{2}\mathbf{I}_p \end{pmatrix}^{-1} \begin{pmatrix} \frac{1}{n}\mathbf{Y} + \frac{r_2}{2}\mathbf{D}^T\mathbf{s} - \frac{1}{2}\mathbf{D}^T\mathbf{q}_2 \\ \frac{1}{n}\mathbf{X}^T\mathbf{Y} + \frac{r_3}{2}\mathbf{w} - \frac{1}{2}\mathbf{q}_3 \end{pmatrix}$$

If the dimension is large, we can use a coordinate descent strategy like before.

$$\boldsymbol{\beta}^{(k+1)} = \left(\frac{1}{n}\mathbf{X}^T\mathbf{X} + \frac{r_3}{2}\mathbf{I}_p \right)^{-1} \left(\frac{1}{n}\mathbf{X}^T\mathbf{Y} - \frac{1}{n}\mathbf{X}^T\boldsymbol{\mu} + \frac{r_3}{2}\mathbf{w} - \frac{1}{2}\mathbf{q}_3 \right)$$

and

$$\boldsymbol{\mu}^{(k+1)} = \left(\frac{1}{n}\mathbf{I}_n + \frac{r_2}{2}\mathbf{D}^T\mathbf{D} \right)^{-1} \left(\frac{1}{n}\mathbf{Y} - \frac{1}{n}\mathbf{X}\boldsymbol{\beta} + \frac{r_2}{2}\mathbf{D}^T\mathbf{s} - \frac{1}{2}\mathbf{D}^T\mathbf{q}_2 \right)$$

A more accurate solution, which is essentially equal to solve the whole linear system would be **to be add**.

A Re-deduction of the Algorithm with More Parameters

In this section, we again give the details of the algorithm, but with more flexibility in the loss function. We consider the loss function

$$L(\boldsymbol{\mu}, \boldsymbol{\beta}) = \frac{1}{a} \sum_{i=1}^n \rho(y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta}) + b \sum_{i < j} P_{\lambda_1}(\mu_i - \mu_j) + c \sum_{j=1}^p P_{\lambda_2}(\beta_j)$$

where a , b and c are constants. The augmented lagrangian format is

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\beta}, \mathbf{z}, \mathbf{s}, \mathbf{w}) &= \frac{1}{a} \sum_{i=1}^n \rho(z_i) + b \sum_{i < j} P_{\lambda_1}(s_{ij}) + c \sum_{j=1}^p P_{\lambda_2}(w_j) \\ &\quad + \frac{r_1}{2} \|\mathbf{y} - \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}\|_2^2 + \frac{r_2}{2} \|\mathbf{D}\boldsymbol{\mu} - \mathbf{s}\|_2^2 + \frac{r_3}{2} \|\boldsymbol{\beta} - \mathbf{w}\|_2^2 \\ &\quad + \langle \mathbf{y} - \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}, \mathbf{q}_1 \rangle + \langle \mathbf{D}\boldsymbol{\mu} - \mathbf{s}, \mathbf{q}_2 \rangle + \langle \boldsymbol{\beta} - \mathbf{w}, \mathbf{q}_3 \rangle \end{aligned}$$

The update of $\boldsymbol{\beta}$ and $\boldsymbol{\mu}$ is unchanged as before. But the update of \mathbf{z} , \mathbf{s} and \mathbf{w} will be affected by the choice of a , b and c .

A.1 Update steps for $\boldsymbol{\beta}$ and $\boldsymbol{\mu}$

We can update $\boldsymbol{\beta}$ and $\boldsymbol{\mu}$ in a coordinate descent fashion. **NOTE:** the previous updating steps of $\boldsymbol{\beta}$ and $\boldsymbol{\mu}$ is a special case of coordinate-descent algorithm, with max number of iteration set fixed to 1.

A.1.1 Update $\boldsymbol{\beta}$, Coordinate Descent

If $p \leq n$, then

$$\boldsymbol{\beta}^{(k+1)} = (r_1 \mathbf{X}^T \mathbf{X} + r_3 \mathbf{I}_p)^{-1} (r_1 \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu} - \mathbf{z}) + r_3 \mathbf{w} + \mathbf{X}^T \mathbf{q}_1 - \mathbf{q}_3).$$

If $p > n$, then

$$\boldsymbol{\beta}^{(k+1)} = \frac{1}{r_3} \left[\mathbf{I}_p - r_1 \mathbf{X}^T (r_1 \mathbf{X} \mathbf{X}^T + r_3 \mathbf{I}_n)^{-1} \mathbf{X} \right] (r_1 \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu} - \mathbf{z}) + r_3 \mathbf{w} + \mathbf{X}^T \mathbf{q}_1 - \mathbf{q}_3).$$

A.1.2 Update $\boldsymbol{\mu}$, Coordinate Descent

The update of $\boldsymbol{\mu}$ is simply

$$\boldsymbol{\mu} = (r_1 \mathbf{I}_n + r_2 \mathbf{D}^T \mathbf{D})^{-1} (r_1 (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}) + r_2 \mathbf{D}^T \mathbf{s} + \mathbf{q}_1 - \mathbf{D}^T \mathbf{q}_2).$$

A.1.3 Update β and μ simultaneously

It means to minimize the following equation with respect to β and μ :

$$\begin{aligned}
& \frac{r_1}{2} (\mathbf{y} - \mu - \mathbf{X}\beta)^T (\mathbf{y} - \mu - \mathbf{X}\beta) + \frac{r_2}{2} (\mathbf{D}\mu - \mathbf{s})^T (\mathbf{D}\mu - \mathbf{s}) + \frac{r_3}{2} (\beta - \mathbf{w})^T (\beta - \mathbf{w}) \\
& + (\mathbf{y} - \mu - \mathbf{X}\beta)^T \mathbf{q}_1 + (\mathbf{D}\mu - \mathbf{s})^T \mathbf{q}_2 + (\beta - \mathbf{w})^T \mathbf{q}_3 \\
& \propto \frac{r_1}{2} (\mu^T \mu + \beta^T \mathbf{X}^T \mathbf{X} \beta - 2\mu^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + 2\beta^T \mathbf{X}^T \mu + 2\mu^T \mathbf{z} + 2\beta^T \mathbf{X}^T \mathbf{z}) \\
& + \frac{r_2}{2} (\mu^T \mathbf{D}^T \mathbf{D} \mu - 2\mu^T \mathbf{D}^T \mathbf{s}) + \frac{r_3}{2} (\beta^T \beta - 2\beta^T \mathbf{w}) - \mu^T \mathbf{q}_1 - \beta^T \mathbf{X}^T \mathbf{q}_1 + \mu^T \mathbf{D}^T \mathbf{q}_2 + \beta^T \mathbf{q}_3 \\
& = (\mu^T \quad \beta^T) \begin{pmatrix} \frac{r_1}{2} \mathbf{I}_n & \frac{r_1}{2} \mathbf{X} \\ \frac{r_1}{2} \mathbf{X}^T & \frac{r_1}{2} \mathbf{X}^T \mathbf{X} \end{pmatrix} \begin{pmatrix} \mu \\ \beta \end{pmatrix} - r_1 (\mathbf{y} - \mathbf{z})^T (\mathbf{I}_n \quad \mathbf{X}) \begin{pmatrix} \mu \\ \beta \end{pmatrix} \\
& + (\mu^T \quad \beta^T) \begin{pmatrix} \frac{r_2}{2} \mathbf{D}^T \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \frac{r_3}{2} \mathbf{I}_p \end{pmatrix} \begin{pmatrix} \mu \\ \beta \end{pmatrix} - (r_2 \mathbf{s}^T \mathbf{D} + \mathbf{q}_1^T - \mathbf{q}_2^T \mathbf{D} \quad r_3 \mathbf{w}^T + \mathbf{q}_1^T \mathbf{X} - \mathbf{q}_3^T) \begin{pmatrix} \mu \\ \beta \end{pmatrix} \\
& = (\mu^T \quad \beta^T) \begin{pmatrix} \frac{r_1}{2} \mathbf{I}_n + \frac{r_2}{2} \mathbf{D}^T \mathbf{D} & \frac{r_1}{2} \mathbf{X} \\ \frac{r_1}{2} \mathbf{X}^T & \frac{r_1}{2} \mathbf{X}^T \mathbf{X} + \frac{r_3}{2} \mathbf{I}_p \end{pmatrix} \begin{pmatrix} \mu \\ \beta \end{pmatrix} \\
& - (r_1 (\mathbf{y} - \mathbf{z})^T + r_2 \mathbf{s}^T \mathbf{D} + \mathbf{q}_1^T - \mathbf{q}_2^T \mathbf{D}, \quad r_1 (\mathbf{y} - \mathbf{z})^T \mathbf{X} + r_3 \mathbf{w}^T + \mathbf{q}_1^T \mathbf{X} - \mathbf{q}_3^T) \begin{pmatrix} \mu \\ \beta \end{pmatrix}
\end{aligned}$$

And the update is

$$\begin{pmatrix} \mu \\ \beta \end{pmatrix} = \begin{pmatrix} r_1 \mathbf{I}_n + r_2 \mathbf{D}^T \mathbf{D} & r_1 \mathbf{X} \\ r_1 \mathbf{X}^T & r_1 \mathbf{X}^T \mathbf{X} + r_3 \mathbf{I}_p \end{pmatrix}^{-1} \begin{pmatrix} r_1 (\mathbf{y} - \mathbf{z}) + r_2 \mathbf{D}^T \mathbf{s} + \mathbf{q}_1 - \mathbf{D}^T \mathbf{q}_2 \\ r_1 \mathbf{X}^T (\mathbf{y} - \mathbf{z}) + r_3 \mathbf{w} + \mathbf{X}^T \mathbf{q}_1 - \mathbf{q}_3 \end{pmatrix}$$

A.2 Update of \mathbf{z}

The update of \mathbf{z} takes the same form as before, except for replacing n with a . Hence for L_1 loss, it's

$$z_i^{(k+1)} = S \left(y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1}, \frac{1}{ar_1} \right)$$

For Huber loss with parameter c_h , the update is

$$z_i^{(k+1)} = \begin{cases} y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1} - \frac{c_h}{ar_1} & \frac{c_h}{ar_1} + c_h < y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1} \\ \frac{y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1}}{1 + \frac{1}{ar_1}} & -\frac{c_h}{ar_1} - c_h \leq y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1} \leq \frac{c_h}{ar_1} + c_h \\ y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1} + \frac{c_h}{ar_1} & y_i - \mu_i - \mathbf{x}_i^T \beta + \frac{q_{1,i}}{r_1} < -\frac{c_h}{ar_1} - c_h \end{cases}$$

A.3 Update of \mathbf{s}

For Lasso penalty with parameter λ_1 , the update is given by

$$s_{ij}^{(k+1)} = S \left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{b\lambda_1}{r_2} \right).$$

For SCAD penalty with paramters λ_1 and γ_1 , the update is given by

$$s_{ij}^{(k+1)} = \begin{cases} S\left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{b\lambda_1}{r_2}\right) & \left|\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}\right| \leq \left(1 + \frac{b}{r_2}\right)\lambda_1 \\ \frac{S\left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{b\gamma_1\lambda_1}{r_2(\gamma_1-1)}\right)}{1 - \frac{b}{r_2(\gamma_1-1)}} & \left(1 + \frac{b}{r_2}\right)\lambda_1 < \left|\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}\right| \leq \gamma_1\lambda_1, \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \left|\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}\right| > \gamma_1\lambda_1 \end{cases}$$

and the additional condition to guarantee the convergence during the algorithm is $r_2 > \frac{b}{\gamma_1-1}$.

For MCP penalty with parameters λ_1 and γ_1 , the update is

$$s_{ij}^{(k+1)} = \begin{cases} \frac{S\left(\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}, \frac{b\lambda_1}{r_2}\right)}{1 - \frac{b}{r_2\gamma_1}} & \left|\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}\right| \leq \gamma_1\lambda_1, \\ \mu_i - \mu_j + \frac{q_{2,ij}}{r_2} & \left|\mu_i - \mu_j + \frac{q_{2,ij}}{r_2}\right| > \gamma_1\lambda_1 \end{cases},$$

and the additional condition is $r_2 > \frac{b}{\gamma_1}$.

A.4 Update of w

For Lasso penalty with paramter λ_2 , the update is giiven by

$$w_j^{(k+1)} = S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{\lambda_2}{r_3}\right).$$

For SCAD penalty with parameters λ_2 and γ_2 , the update is given by

$$w_j^{(k+1)} = \begin{cases} S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{c\lambda_2}{r_3}\right) & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| \leq \left(1 + \frac{c}{r_3}\right)\lambda_2 \\ \frac{S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{c\gamma_2\lambda_2}{r_3(\gamma_2-1)}\right)}{1 - \frac{c}{r_3(\gamma_2-1)}} & \left(1 + \frac{c}{r_3}\right)\lambda_2 < \left|\beta_j + \frac{q_{3,j}}{r_3}\right| \leq \gamma_2\lambda_2, \\ \beta_j + \frac{q_{3,j}}{r_3} & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| > \gamma_2\lambda_2 \end{cases}$$

and the additional condition is $r_2 > \frac{c}{\gamma_2-1}$.

For MCP penalty with parameters λ_2 and γ_2 , the update is given by

$$w_j^{(k+1)} = \begin{cases} \frac{S\left(\beta_j + \frac{q_{3,j}}{r_3}, \frac{c\lambda_2}{r_3}\right)}{1 - \frac{c}{r_3\gamma_2}} & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| \leq \gamma_2\lambda_2, \\ \beta_j + \frac{q_{3,j}}{r_3} & \left|\beta_j + \frac{q_{3,j}}{r_3}\right| > \gamma_2\lambda_2 \end{cases},$$

and the additional condition is $r_3 > \frac{c}{\gamma_2}$.

A.5 Choice of $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$

The initial values of $\lambda_1^{(0)}$ and $\lambda_2^{(0)}$ which shrinks all penalty values to 0 is also affected by the choice of a , b and c , but the deduction is quite similar to before. In summary

$$\lambda_2^{(0)} \geq \frac{1}{ac} \left\| \sum_{i=1}^n \psi(y_i - m) \mathbf{x}_i \right\|_{\infty},$$

and

$$\lambda_1^{(0)} \geq \frac{1}{ab} \left\| \mathbf{D} (\mathbf{D}^T \mathbf{D})^{-1} \begin{pmatrix} \psi(y_1 - m) \\ \vdots \\ \psi(y_n - m) \end{pmatrix} \right\|_{\infty},$$

where $(\mathbf{D}^T \mathbf{D})^{-1}$ is the generalized inverse, $\psi = \partial \rho$ and m is the estimating value of μ when all penalties are pushed to 0, i.e.

$$m = \underset{\mu}{\operatorname{argmin}} \sum_{i=1}^n \rho(y_i - \mu).$$

To simplify the computation, we can choose

$$\lambda_1^{(0)} \geq \frac{2}{abn} \left\| \begin{pmatrix} \psi(y_1 - m) \\ \vdots \\ \psi(y_n - m) \end{pmatrix} \right\|_{\infty}$$