

ADMM for High-Dimensional Sparse Penalized Quantile Regression

Qinqin Zou and Qian Tang

University of Iowa

May 8, 2023

Quantile Regression

- τ th quantile: $\tau \in (0, 1)$, $Q_Y(\tau) = \inf_y \{y : F(y) \geq \tau\}$.
- **Quantile Regression** (Koenker and Bassett, 1978):

$$\min_{b, \beta} \sum_{i=1}^n \rho_{\tau}(y_i - b - x_i' \beta); \text{ check loss: } \rho_{\tau}(t) = t(\tau - I(t < 0)).$$

- Median regression ($\tau = 0.5$):

$$\min_{b, \beta} \sum_{i=1}^n |y_i - b - x_i' \beta|$$

Sparse Penalized Quantile Regression

- Sparse penalized quantile regression:

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \rho_{\tau}(y_i - \mathbf{x}_i^T \beta) + \sum_{j=1}^n p_{\lambda}(|\beta_j|)$$

- $p_{\lambda}(\cdot)$, $\lambda > 0$ is the penalty function: lasso penalty
- Weighted L_1 -penalized quantile regression:

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \rho_{\tau}(y_i - \mathbf{x}_i^T \beta) + \lambda \|\mathbf{w} \circ \beta\|_1$$

- lasso penalized quantile regression: $\mathbf{w} = \mathbf{1}_p$
- adaptive lasso penalized quantile regression: $w_j = \left(\left| \hat{\beta}_j^{\text{lasso}} \right| + 1/n \right)^{-1}$,
 $j = 1, \dots, p$

Weighted L_1 -penalized quantile regression

- Weighted L_1 -penalized quantile regression:

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \rho_{\tau}(y_i - x_i^T \beta) + \lambda \|\mathbf{w} \circ \beta\|_1 \quad (1)$$

- problem (1) is equivalent to:

$$\begin{aligned} & \min_{\beta, \mathbf{z}} \mathbb{Q}_{\tau}(\mathbf{z}) + \lambda \|\mathbf{w} \circ \beta\|_1 \\ & \text{subject to } \mathbf{X}\beta + \mathbf{z} = \mathbf{y}. \end{aligned} \quad (2)$$

where $\mathbb{Q}_{\tau}(\mathbf{z}) = (1/n) \sum_{i=1}^n \rho_{\tau}(z_i)$ and $\mathbf{z} = \mathbf{y} - \mathbf{X}\beta$

ADMM Algorithm

Fix $\sigma > 0$ and the augmented Lagrangian function of (2) is:

$$\begin{aligned}\mathcal{L}_\sigma(\boldsymbol{\beta}, \mathbf{z}, \boldsymbol{\theta}) := & \mathbb{Q}_\tau(\mathbf{z}) + \lambda \|\mathbf{w} \circ \boldsymbol{\beta}\|_1 - \langle \boldsymbol{\theta}, \mathbf{X}\boldsymbol{\beta} \\ & + \mathbf{z} - \mathbf{y} \rangle + \frac{\sigma}{2} \|\mathbf{X}\boldsymbol{\beta} + \mathbf{z} - \mathbf{y}\|_2^2\end{aligned}$$

Following Boyd et al. (2011), the iterations for the standard ADMM algorithm are given by

$$\boldsymbol{\beta} \text{ step : } \boldsymbol{\beta}^{k+1} := \arg \min_{\boldsymbol{\beta}} \lambda \|\mathbf{w} \circ \boldsymbol{\beta}\|_1 - \langle \boldsymbol{\theta}^k, \mathbf{X}\boldsymbol{\beta} \rangle + \frac{\sigma}{2} \|\mathbf{X}\boldsymbol{\beta} + \mathbf{z}^k - \mathbf{y}\|_2^2$$

$$\mathbf{z} \text{ step : } \mathbf{z}^{k+1} := \arg \min \mathbb{Q}_\tau(\mathbf{z}) - \langle \boldsymbol{\theta}^k, \mathbf{z} \rangle + \frac{\sigma}{2} \|\mathbf{z} + \mathbf{X}\boldsymbol{\beta}^{k+1} - \mathbf{y}\|_2^2$$

$$\boldsymbol{\theta} \text{ step : } \boldsymbol{\theta}^{k+1} := \boldsymbol{\theta}^k - \sigma (\mathbf{X}\boldsymbol{\beta}^{k+1} + \mathbf{z}^{k+1} - \mathbf{y}).$$

ADMM Algorithm

- the proximal mapping: $\text{Prox}_{\rho_\tau}[\xi, \alpha] := \arg \min_{u \in \mathbb{R}} \rho_\tau(u) + \frac{\alpha}{2}(u - \xi)^2$
- For $i = 1, \dots, n$, we have

$$\begin{aligned} z_i^{k+1} &:= \arg \min_{z_i} \frac{1}{n} \rho_\tau(z_i) - \theta_i^k z_i + \frac{\sigma}{2} \left(z_i + \mathbf{x}_i^T \boldsymbol{\beta}^{k+1} - y_i \right)^2 \\ &= \text{Prox}_{\rho_\tau} \left[y_i - \mathbf{x}_i^T \boldsymbol{\beta} + \frac{1}{\sigma} \theta_i^k, n\sigma \right] \end{aligned} \quad (3)$$

Lemma 1

Given $\tau \in (0, 1)$ and $\alpha > 1$, the proximal mapping

$\text{Prox}_{\rho_\tau}[\xi, \alpha] := \arg \min_{u \in \mathbb{R}} \rho_\tau(u) + \frac{\alpha}{2}(u - \xi)^2$ has explicit expression,

$$\text{Prox}_{\rho_\tau}[\xi, \alpha] = \begin{cases} \xi - \frac{\tau}{\alpha}, & \text{if } \xi > \frac{\tau}{\alpha} \\ 0, & \text{if } \frac{\tau-1}{\alpha} \leq \xi \leq \frac{\tau}{\alpha} \\ \xi - \frac{\tau-1}{\alpha}, & \text{if } \xi < \frac{\tau-1}{\alpha}. \end{cases}$$

ADMM Algorithm

- β step : $\beta^{k+1} := \arg \min_{\beta} \lambda \|\mathbf{w} \circ \beta\|_1 - \langle \boldsymbol{\theta}^k, \mathbf{X}\beta \rangle + \frac{\sigma}{2} \|\mathbf{X}\beta + \mathbf{z}^k - \mathbf{y}\|_2^2$
- Augmented β step;

$$\begin{aligned} \beta^{k+1} := & \arg \min_{\beta} \lambda \|\mathbf{w} \circ \beta\|_1 - \langle \boldsymbol{\theta}^k, \mathbf{X}\beta \rangle \\ & + \frac{\sigma}{2} \|\mathbf{X}\beta + \mathbf{z}^k - \mathbf{y}\|_2^2 \\ & + \frac{1}{2} \|\beta - \beta^k\|_S^2 \end{aligned}$$

- $\|\mathbf{v}\|_S^2 := \langle \mathbf{v}, \mathbf{S}\mathbf{v} \rangle$
- $\mathbf{S} = \sigma (\eta \mathbf{I}_p - \mathbf{X}^\top \mathbf{X})$ with $\eta \geq \Lambda_{\max}(\mathbf{X}^\top \mathbf{X})$

Soft-thresholding

For the problem,

$$\min_{\mathbf{X}} \|\mathbf{X} - B\|_2^2 + \lambda \|\mathbf{X}\|_1$$

the solution can be written as

$$\hat{x}_i = \text{shrink}[b_i, \lambda/2] = \text{sgn}(b_i) \max(|b_i| - \lambda/2, 0).$$

$$\begin{aligned} \beta^{k+1} &= \arg \min_{\beta} \lambda \|\mathbf{w} \circ \beta\|_1 \\ &\quad + \frac{\sigma\eta}{2} \left\| \beta - \frac{\sigma\eta\beta^k + \mathbf{X}^T (\boldsymbol{\theta}^k + \sigma y - \sigma\mathbf{X}\beta^k - \sigma\mathbf{z}^k)}{\sigma\eta} \right\|_2^2 \\ &= \left(\text{Shrink} \left[\beta_j^k + \frac{1}{\sigma\eta} \mathbf{X}_j^T (\boldsymbol{\theta}^k + \sigma y - \sigma\mathbf{X}\beta^k - \sigma\mathbf{z}^k), \frac{\lambda w_j}{\sigma\eta} \right] \right)_{1 \leq j \leq p} \end{aligned} \quad (4)$$

Algorithm 1 pADMM – Proximal ADMM algorithm for solving the weighted L_1 -penalized quantile regression.

1. Initialize the algorithm with $(\boldsymbol{\beta}^0, \mathbf{z}^0, \boldsymbol{\theta}^0)$.
2. For $k = 0, 1, 2, \dots$, repeat steps (2.1)–(2.3) until the convergence criterion is met.

(2.1) Update $\boldsymbol{\beta}^{k+1} \leftarrow$

$$\left(\text{Shrink} \left[\boldsymbol{\beta}_j^k + \frac{1}{\sigma\eta} \mathbf{X}_j^T (\boldsymbol{\theta}^k + \sigma \mathbf{y} - \sigma \mathbf{X} \boldsymbol{\beta}^k - \sigma \mathbf{z}^k), \frac{\lambda w_j}{\sigma\eta} \right] \right)_{1 \leq j \leq p}.$$

(2.2) Update $\mathbf{z}^{k+1} \leftarrow \left(\text{Prox}_{\rho\tau} [y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{k+1} + \sigma^{-1} \theta_i^k, n\sigma] \right)_{1 \leq i \leq n}.$

(2.3) Update $\boldsymbol{\theta}^{k+1} \leftarrow \boldsymbol{\theta}^k - \gamma\sigma (\mathbf{X} \boldsymbol{\beta}^{k+1} + \mathbf{z}^{k+1} - \mathbf{y})$.

- $\eta \geq \Lambda_{\max}(\mathbf{X}^T \mathbf{X})$
- $\boldsymbol{\beta}$ step : $\boldsymbol{\beta}^{k+1} := \arg \min_{\boldsymbol{\beta}} \lambda \|\mathbf{w} \circ \boldsymbol{\beta}\|_1 - \langle \boldsymbol{\theta}^k, \mathbf{X} \boldsymbol{\beta} \rangle + \frac{\sigma}{2} \|\mathbf{X} \boldsymbol{\beta} + \mathbf{z}^k - \mathbf{y}\|_2^2$

Algorithm 2 scdADMM – Sparse coordinate descent ADMM algorithm for solving the weighted L_1 -penalized quantile regression with coordinate descent steps.

1. Initialize the algorithm with $(\boldsymbol{\beta}^0, \mathbf{z}^0, \boldsymbol{\theta}^0)$.
2. For $k = 0, 1, 2, \dots$, repeat Steps (2.1)–(2.3) until the convergence criterion is met.

(2.1) Carry out the coordinate descent Steps (2.1.1)–(2.1.3).

(2.1.1) Initialize $\boldsymbol{\beta}^{k,0} = \boldsymbol{\beta}^k$.

(2.1.2) For $m = 0, 1, 2, \dots$, repeat Step (2.1.2.1) until convergence.

(2.1.2.1) For $j = 1, \dots, p$, update

$$\beta_j^{k,m+1} \leftarrow \frac{\text{Shrink} \left[\sum_{i=1}^n x_{ij} \left\{ \theta_i^k + \sigma \left(y_i - z_i^k - \sum_{t \neq j} x_{it} \beta_t^{k,m+1(t < j)} \right) \right\}, \lambda w_j \right]}{\sigma \|X_j\|_2^2}.$$

(2.1.3) Set $\boldsymbol{\beta}^{k+1} \leftarrow \boldsymbol{\beta}^{k,m+1}$.

(2.2) Update $\mathbf{z}^{k+1} \leftarrow \left(\text{Prox}_{\rho\tau} [y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{k+1} + \sigma^{-1} \theta_i^k, n\sigma] \right)_{1 \leq i \leq n}$.

(2.3) Update $\boldsymbol{\theta}^{k+1} \leftarrow \boldsymbol{\theta}^k - \sigma (\mathbf{X} \boldsymbol{\beta}^{k+1} + \mathbf{z}^{k+1} - \mathbf{y})$.

Theorem 1

For given $\lambda > 0, \sigma > 0, 0 < \tau < 1, 0 < \gamma < (\sqrt{5} + 1)/2$ and a component-wisely nonnegative weight vector w , let $\{(\beta^k, z^k, \theta^k)\}$ be generated by the pADMM algorithm as described in Algorithm 1. Then, the sequence $\{(\beta^k, z^k), k = 0, 1, 2, \dots\}$ converges to an optimal solution (β^*, z^*) to (2) and $\{\theta^k, k = 0, 1, 2, \dots\}$ converges to an optimal solution θ^* to the dual problem of (2). Equivalently, $\{\beta^k, k = 0, 1, 2, \dots\}$ converges to a global minimizer of problem (1).

Note that the convergence of the scdADMM algorithm (Algorithm 2) can be directly obtained from Boyd et al. (2011).

Implementation Details

- R package: FHDQR

- λ sequence: $\lambda_k = \lambda_{\max}^{\frac{M-k}{M-1}} \lambda_{\min}^{\frac{k-1}{M-1}}, k = 1, \dots, M, \lambda_{\min} = \delta \lambda_{\max}$

- Implementation skill: warm-start

- Stop criterion:

$$\left\| X\beta^k + z^k - y \right\|_2 \leq \sqrt{n}\epsilon_1 + \epsilon_2 \max \left\{ \left\| X\beta^k \right\|_2, \left\| z^k \right\|_2, \|y\|_2 \right\},$$
$$\sigma \left\| X^T \left(z^k - z^{k-1} \right) \right\|_2 \leq \sqrt{p}\epsilon_1 + \epsilon_2 \left\| X^T \theta^k \right\|_2,$$

Data generating models

- Example 1:

$$Y = \sum_{j=1}^p X_j \beta_j + k \cdot \varepsilon$$

- $(X_1, \dots, X_p)^T \sim N(0, \Sigma)$ with $\Sigma = (\alpha + (1 - \alpha)I(i = j))_{p \times p}$
- $\alpha \in \{0, 0.1, 0.2, 0.5, 0.9, 0.95\}$
- $\beta_j = (-1)^j \exp(-(2j - 1)/20)$
- $\varepsilon \sim N(0, 1)$
- $k = 3$
- $n = 100, p = 1000$

- Example 2:

$$y_i = x_i^T \beta^* + \varepsilon_i, \quad x_i \sim N(0, \Sigma_x), \quad i = 1, \dots, n$$

- $\beta^* = (2, 0, 1.5, 0, 0.8, 0, 0, 1, 0, 1.75, 0, 0, 0.75, 0, 0, 0.3, 0_{p-16}^T)^T$
- $n=200, p=1000$

- Timing Comparisons:
 - Penalty function: lasso
 - Solvers: FHDQR, hqreg, quantreg
- Finite-Sample Performance:
 - Penalized quantile regression vs. penalized least squares regression
 - Penalty function: lasso, adaptive lasso and SCAD

Example 1: Computation time

	Correlation (α)					
	0.00	0.10	0.20	0.50	0.90	0.95
$\tau = 0.25$						
quantreg	300+	300+	300+	300+	300+	300+
hqreg	9.52	9.32	9.82	9.86	7.05	5.63
pADMM	0.57	4.38	5.85	12.14	18.43	11.62
scdADMM	1.41	1.37	1.36	1.23	1.00	0.94
$\tau = 0.50$						
quantreg	300+	300+	300+	300+	300+	300+
hqreg	6.88	6.64	6.89	7.92	8.65	5.29
pADMM	0.62	5.36	7.85	15.47	30.34	21.66
scdADMM	1.26	1.20	1.26	1.18	1.19	0.91
$\tau = 0.75$						
quantreg	300+	300+	300+	300+	300+	300+
hqreg	8.65	8.34	8.81	8.87	8.37	5.76
pADMM	0.55	4.45	6.26	12.12	21.49	16.40
scdADMM	1.42	1.39	1.48	1.34	1.15	1.20

- Different algorithms all yield the same (numerically speaking) objective function values.

Example 2: Computation time($\tau = 0.5$)

Covariance	Method	Error			
		$N(0, 2)$	Laplace	$\sqrt{2} \times t_4$	Cauchy
$\Sigma_x = I$	quantreg	400+	400+	400+	400+
	hqreg	10.45	10.45	13.95	32.96
	scdADMM	3.03	2.88	3.39	5.88
	pADMM	1.52	1.47	1.46	0.46
$AR_{0.5}$	quantreg	400+	400+	400+	400+
	hareg	11.19	10.48	14.11	24.56
	scdADMM	3.76	3.47	4.11	5.85
	pADMM	1.83	1.76	1.77	0.55
$CS_{0.5}$	quantreg	400+	400+	400+	400+
	hqreg	13.70	10.11	14.32	15.04
	scdADMM	7.11	6.77	7.68	8.49
	pADMM	19.91	17.96	20.12	5.65

- $AR_{0.5} : \Sigma_x = (0.5^{|i-j|})$
- $CS_{0.5} : \Sigma_x = (0.5 + (1 - 0.5)I(i = j))$

Example 2: Finite-Sample Performance (L_1, L_2 losses)

		$\Sigma = (0.5 + 0.5I(i = j))$		
		Lasso	Alasso	SCAD
$N(0, 2)$	LS	4.229, 0.959	1.717, 0.610	1.286, 0.570
	QR(0.50)	5.744, 1.177	1.710, 0.718	1.681, 0.674
	QR(0.75)	5.630, 1.236	2.016, 0.835	2.170, 0.824
Laplace	LS	4.217, 0.973	1.753, 0.620	1.294, 0.578
	QR(0.50)	4.120, 0.887	1.110, 0.494	0.991, 0.408
	QR(0.75)	5.207, 1.103	1.667, 0.700	1.609, 0.666
$\sqrt{2} \times t_4$	LS	6.027, 1.342	3.706, 1.114	2.265, 0.975
	QR(0.50)	5.941, 1.276	2.045, 0.842	1.900, 0.780
	QR (0.75)	6.485, 1.456	2.819, 1.090	2.959, 1.112
Cauchy	LS	20.931, 4.650	32.018, 8.487	472.891, 46.890
	QR(0.50)	5.905, 1.324	2.497, 0.978	2.257, 0.868
	QR(0.75)	8.193, 1.904	4.563, 1.598	4.904, 1.684

A Real Data Example

- Goal: study how the expression of *TRIM32* depends on the expressions of 3000 genes.
- $n=120$, $p=3000$
- Timing results:

τ	0.25	0.50	0.75
quantreg	5000+	5000+	5000+
hqreg	4.97	4.09	4.56
pADMM	351.93	401.76	347.89
scdADMM	1.68	1.15	1.09

A Real Data Example

- Prediction error $= (1/40) \sum_{i \in \text{validation}} \rho_{\tau} \left(y_i - \hat{\beta}_0 - \mathbf{x}_i^T \hat{\boldsymbol{\beta}} \right)$
- 50 random partitions
- All the results are averaged over 50 runs for the random partition columns.

τ	All data #genes	Random partition	
		Ave. #genes	Prediction error
0.25	14	15.00(1.26)	0.0351(0.0014)
0.50	23	24.16(2.38)	0.0395(0.0010)
0.75	14	11.22(1.07)	0.0671(0.0196)

Thank you