

# Coordinate Descent Algorithm for Least Absolute Deviations Regression

Zehaan Naik

August 25, 2025

## 1 Introduction

Linear regression is a foundational statistical method for modeling the relationship between a dependent variable and one or more explanatory variables. The standard approach for estimating the model parameters, Ordinary Least Squares (OLS), is highly effective under the assumption of normally distributed errors. OLS provides the best linear unbiased estimates by minimizing the sum of the squared residuals. This objective function is mathematically convenient, yielding a well-known analytical solution:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

However, the reliance on squared errors makes OLS highly sensitive to outliers. A single data point far from the general trend can disproportionately influence the regression line, leading to a model that poorly represents the bulk of the data. This lack of robustness is a significant drawback in many real-world applications where datasets are often contaminated with anomalous observations. To address this, robust regression methods have been developed.

Among the most prominent is the Least Absolute Deviations (LAD) method, also known as  $L_1$  regression. Instead of minimizing the sum of squared errors, LAD regression seeks to minimize the sum of the absolute errors:

$$\min_{\beta} \sum_{i=1}^n |y_i - \mathbf{x}_i^T \beta|$$

This objective function is inherently more robust to outliers, as the influence of a large residual is linear rather than quadratic. Despite its desirable properties, the LAD objective function is non-differentiable at the origin, which means there is no simple analytical solution akin to the one for OLS. Historically, the LAD estimation problem has been solved using linear programming techniques. The problem can be reformulated as a linear program and solved using standard algorithms like the simplex method Charnes et al. (1955); Wagner (1959); Barrodale and Roberts (1973, 1974); Koenker and Bassett (1978).

While these methods provide exact solutions, they can be computationally intensive and complex to implement from scratch, often requiring specialized optimization software. This paper introduces an iterative algorithm for computing LAD estimates that is both computationally efficient and remarkably simple to implement. Our proposed method is a form of coordinate descent (Koenker and d’Orey, 1987; Peng and Wang, 2015; Mkhadri et al., 2017) that iteratively solves a series of one-dimensional minimization problems, each of which has a straightforward solution based on the weighted median. We demonstrate that this approach is a valid descent algorithm, guaranteeing convergence to the optimal solution. The primary advantage of our algorithm lies in its simplicity and speed, making robust LAD regression more accessible without the overhead of complex linear programming solvers.

## 2 Background

The proposed algorithm leverages the principle of coordinate descent, breaking the multi-dimensional LAD optimization problem into a sequence of simpler, one-dimensional problems. The solvability of our algorithm hinges on the fact that these one-dimensional sub-problems have well-known, closed-form solutions. In this section, we establish the theoretical foundation by detailing the solutions to these core optimization tasks.

### 2.1 Minimizing Deviations from a Constant

The first and most fundamental sub-problem is to find a single constant,  $a$ , that best represents a set of observations  $\{y_1, \dots, y_n\}$  in the  $L_1$  sense. This requires minimizing the objective function:

$$L(a) = \sum_{i=1}^n |y_i - a| \tag{1}$$

It is a classical result in statistics that the minimizer of this function is the **median** of the observations  $\{y_i\}$ . To show this, we can examine the subgradient of the loss function. The derivative of  $|y_i - a|$  with respect to  $a$  is  $-\text{sgn}(y_i - a)$ , where  $\text{sgn}(\cdot)$  is the sign function. The subgradient of  $L(a)$  is therefore:

$$\partial L(a) = \sum_{i=1}^n -\text{sgn}(y_i - a) \quad (2)$$

The minimum is achieved when  $0 \in \partial L(a)$ , which implies that the number of observations greater than  $a$  must equal the number of observations less than  $a$ . This is precisely the definition of the sample median (Koenker and Bassett, 1978).

## 2.2 Minimizing Deviations for a No-Intercept Model

The second sub-problem involves finding the optimal slope coefficient,  $b$ , for a simple regression model that passes through the origin. The objective is to minimize:

$$L(b) = \sum_{i=1}^n |y_i - bx_i| \quad (3)$$

This problem is slightly more complex than the first, but it can be reframed into a familiar form. By factoring out  $|x_i|$  (assuming  $x_i \neq 0$ ), we can rewrite the objective function as:

$$L(b) = \sum_{i=1}^n |x_i| \left| \frac{y_i}{x_i} - b \right| \quad (4)$$

This reveals that the problem is equivalent to finding the **weighted median** of the ratios  $\{y_i/x_i\}$ , where the weight for each ratio is given by the corresponding  $|x_i|$ . The weighted median is the value  $b$  that partitions the data such that the sum of the weights of the ratios less than  $b$  is equal to the sum of the weights of the ratios greater than  $b$ . This solution is also a direct consequence of the properties of quantile regression (Koenker and Bassett, 1978). These two fundamental results form the building blocks of our iterative algorithm, allowing each step to be performed efficiently and exactly.

### 3 Proposed Methodology

Our approach to solving the LAD regression problem is based on a coordinate descent framework. Instead of confronting the complex, multi-dimensional objective function directly, we break it down into a sequence of simpler, one-dimensional optimization problems that can be solved exactly and efficiently.

#### 3.1 The Algorithm: An Iterative Approach

Consider the general LAD objective function for a model with  $p$  covariates and an intercept, where  $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_p]^T$  and  $\mathbf{x}_i = [1, x_{i1}, \dots, x_{ip}]^T$ :

$$L(\boldsymbol{\beta}) = \sum_{i=1}^n |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| = \sum_{i=1}^n |y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij}| \quad (5)$$

The core idea of our algorithm is to optimize this function with respect to one parameter at a time, while holding all other parameters fixed at their current values. For example, to find the optimal value for a single coefficient  $\beta_j$ , we can define the partial residuals as  $r_i^{(j)} = y_i - \beta_0 - \sum_{k \neq j} \beta_k x_{ik}$ . The optimization problem for  $\beta_j$  then becomes:

$$\min_{\beta_j} \sum_{i=1}^n |r_i^{(j)} - \beta_j x_{ij}| \quad (6)$$

As established in Section 2.2, the solution to this is the weighted median of the ratios  $\{r_i^{(j)}/x_{ij}\}$  with weights  $\{|x_{ij}|\}$ . Similarly, when solving for the intercept  $\beta_0$ , we define residuals  $r_i^{(0)} = y_i - \sum_{j=1}^p \beta_j x_{ij}$  and solve:

$$\min_{\beta_0} \sum_{i=1}^n |r_i^{(0)} - \beta_0| \quad (7)$$

The solution is simply the median of the residuals  $\{r_i^{(0)}\}$ . By cyclically iterating through all parameters and applying these updates, we progressively minimize the total loss function.

#### 3.2 Proof of Convergence

The iterative nature of our algorithm guarantees that the objective function is non-increasing at every step, ensuring convergence to the global minimum. This property is formalized in

the following theorem.

*Lemma 1* (Convexity of the LAD objective). Let

$$L(\beta) = \sum_{i=1}^n |y_i - x_i^\top \beta|,$$

where  $\beta \in \mathbb{R}^{p+1}$ . Then  $L$  is convex in  $\beta$ .

*Proof.* The map  $t \mapsto |t|$  is convex on  $\mathbb{R}$ . For each  $i$ , the residual

$$r_i(\beta) = y_i - x_i^\top \beta$$

is an affine function of  $\beta$ . The composition of a convex function with an affine function is convex, hence  $\beta \mapsto |r_i(\beta)|$  is convex. Since a nonnegative sum of convex functions is convex, it follows that  $L(\beta) = \sum_{i=1}^n |r_i(\beta)|$  is convex.  $\square$

*Corollary 1* (Convexity of each coordinate subproblem). Fix all coordinates of  $\beta$  except  $\beta_j$ . Define

$$g_j(t) = \sum_{i=1}^n \left| r_i^{(j)} - x_{ij}t \right|, \quad \text{where} \quad r_i^{(j)} = y_i - \beta_0 - \sum_{k \neq j} \beta_k x_{ik}.$$

Then  $g_j(t)$  is convex in  $t$ . Consequently, each coordinate update in Algorithm 1 solves a convex (piecewise-linear) optimization problem, and its minimizer set is a closed interval whose endpoints are weighted medians of  $\{r_i^{(j)}/x_{ij}\}$  with weights  $|x_{ij}|$ .

**Theorem 1.** Let  $L(\beta)$  be the LAD objective function. Let  $\beta^k$  be the vector of parameter estimates at the beginning of iteration  $k$ . Let  $\beta^{k+1}$  be the vector of estimates after one full cycle of coordinate-wise updates. Then, the sequence of losses is non-increasing:

$$L(\beta^{k+1}) \leq L(\beta^k)$$

*Proof.* A full iteration  $k$  consists of  $p+1$  updates, one for each parameter from  $\beta_0$  to  $\beta_p$ . Let  $\beta^{k,j}$  denote the parameter vector after updating the  $j$ -th parameter within the  $k$ -th iteration. The starting point is  $\beta^{k,-1} = \beta^k$ . The update for the intercept  $\beta_0$  is:

$$\beta_0^{k+1} = \arg \min_{\beta_0} L(\beta_0, \beta_1^k, \dots, \beta_p^k)$$

By the definition of the argmin operator, the loss cannot increase:

$$L(\beta_0^{k+1}, \beta_1^k, \dots, \beta_p^k) \leq L(\beta_0^k, \beta_1^k, \dots, \beta_p^k)$$

Next, the update for  $\beta_1$  uses the newly updated intercept:

$$\beta_1^{k+1} = \arg \min_{\beta_1} L(\beta_0^{k+1}, \beta_1, \beta_2^k, \dots, \beta_p^k)$$

Note that Lemma 1 gives us a guarantee about the existence of this solution. This implies:

$$L(\beta_0^{k+1}, \beta_1^{k+1}, \beta_2^k, \dots, \beta_p^k) \leq L(\beta_0^{k+1}, \beta_1^k, \beta_2^k, \dots, \beta_p^k)$$

We can continue this chain of reasoning for all  $p + 1$  parameters. Each step is guaranteed to be non-increasing because it solves its one-dimensional sub-problem optimally. After updating the final parameter,  $\beta_p$ , we have  $\beta^{k+1} = [\beta_0^{k+1}, \dots, \beta_p^{k+1}]$ . The sequence of inequalities yields:

$$L(\beta^{k+1}) \leq \dots \leq L(\beta_0^{k+1}, \beta_1^k, \dots, \beta_p^k) \leq L(\beta^k)$$

Thus, the total loss is non-increasing across any full iteration. Since the LAD objective function is convex (Lemma 1 and Corollary 1) and bounded below by zero, this descent property guarantees that the algorithm will converge to the global minimum.  $\square$

The complete procedure is summarized in Algorithm 1. I also implemented this algorithm for a generated dataset to check the performance and how well the parameters converge. I plot the line of best fit obtained using our algorithm, the MAE on parameters, and the MAE on predictions over iterations in Fig 1.

## 4 Experiments

To validate the theoretical properties and practical utility of our proposed coordinate descent algorithm, we conduct a series of experiments on synthetic data. The primary goal of these experiments is to empirically demonstrate the algorithm’s key feature: its robustness to outliers. We compare our model’s performance against scenarios designed to challenge standard regression techniques and highlight the stability and speed of our method.

---

**Algorithm 1** Coordinate Descent for LAD Regression

---

**Initialize:** parameters  $\beta^{(0)} = [\beta_0^{(0)}, \dots, \beta_p^{(0)}]$ , e.g., with zeros or OLS estimates.  
**Set** tolerance  $\epsilon > 0$  and max iterations  $M$ .  
**for**  $k = 1$  to  $M$  **do**  
    Store old parameters:  $\beta^{\text{old}} \leftarrow \beta^{(k-1)}$   
    *// Update intercept*  
    Compute residuals:  $r_i \leftarrow y_i - \sum_{j=1}^p \beta_j^{(k-1)} x_{ij}$  for  $i = 1, \dots, n$ .  
    Update intercept:  $\beta_0^{(k)} \leftarrow \text{median}(r_1, \dots, r_n)$ .  
    *// Update slope coefficients*  
    **for**  $j = 1$  to  $p$  **do**  
        Compute partial residuals:  $r_i^{(j)} \leftarrow y_i - \beta_0^{(k)} - \sum_{l \neq j} \beta_l^{(k-1)} x_{il}$ .  
        Compute ratios:  $z_{ij} \leftarrow r_i^{(j)} / x_{ij}$ .  
        Compute weights:  $w_{ij} \leftarrow |x_{ij}|$ .  
        Update coefficient:  $\beta_j^{(k)} \leftarrow \text{weighted\_median}(\{z_{ij}\}, \{w_{ij}\})$ .  
    **end for**  
    **Check for convergence:**  
    **if**  $\|\beta^{(k)} - \beta^{\text{old}}\|_1 < \epsilon$  **then**  
        **break**  
    **end if**  
**end for**  
**Return**  $\beta^{(k)}$

---

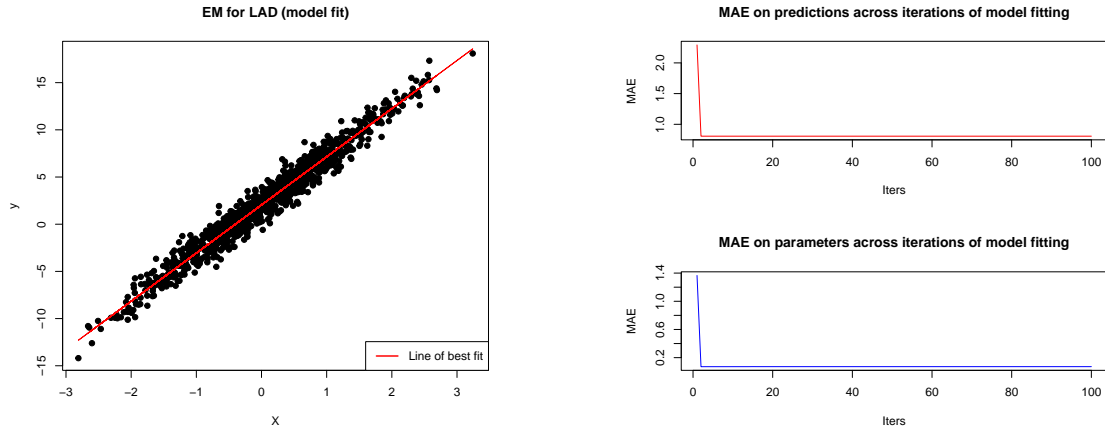


Figure 1: Performance of the coordinate descent algorithm on the generated data set

### 4.1 Performance on Data with Outliers

In the first experiment, we test the algorithm’s performance on a dataset deliberately contaminated with outliers. We generated a simple linear dataset and then introduced several points with large, anomalous y-values that do not follow the underlying linear trend. As demonstrated in Figure 2, the presence of these outliers has a minimal impact on the final LAD estimate. The fitted line correctly captures the true relationship of the majority of the data points, ignoring the influence of the anomalous points. Furthermore, the convergence plots show that both the prediction error (MAE) and the parameter estimates stabilize very quickly, achieving a stable fit in just a few iterations. This experiment clearly illustrates the robustness that motivates the use of LAD regression.

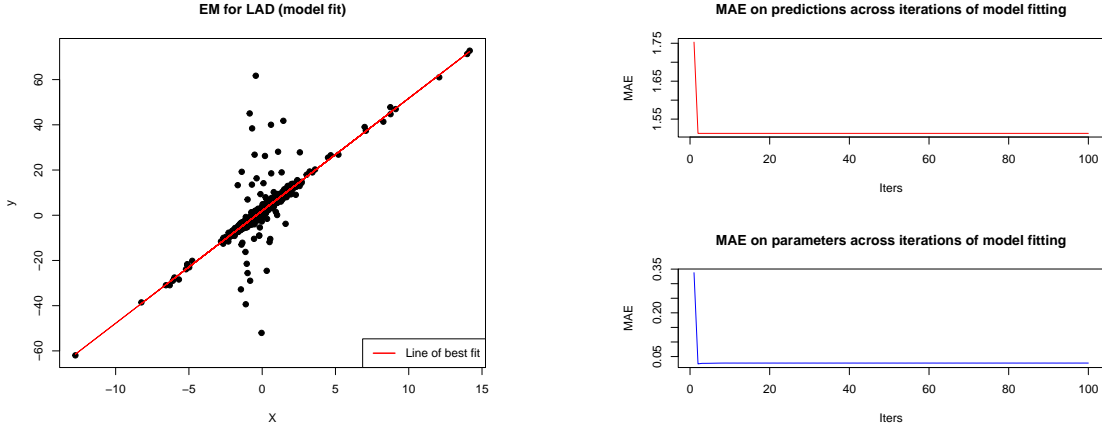


Figure 2: Performance of the coordinate descent algorithm on the generated data set with outliers

### 4.2 Performance in High-Dimensional Settings

To evaluate the scalability and stability of our algorithm, we next consider the challenging high-dimensional scenario where the number of parameters,  $p$ , is close to or exceeds the number of samples,  $n$ . This setting, often referred to as " $p \geq n$ ", can pose significant computational challenges for many statistical methods. We generated datasets with a fixed number of observations ( $n = 1000$ ) and varied the number of parameters across  $p \in \{100, 200, 400, 800, 1200\}$ .

The purpose of this experiment is to empirically demonstrate the algorithm’s robustness



and efficiency as dimensionality increases. The results, shown in Figure 3, are striking. Across all settings where  $p < n$ , the algorithm exhibits extremely rapid convergence, with the Mean Absolute Error (MAE) on the predictions stabilizing in fewer than 20 iterations. This demonstrates that the computational efficiency of our coordinate descent approach does not degrade significantly with the addition of more parameters.

Perhaps the most compelling result is the algorithm’s performance in the ill-conditioned case where  $p > n$ . As shown in Figure 4 for the  $p = 1200$  scenario, the algorithm remains perfectly stable. It converges just as quickly as in the lower-dimensional cases, finding a stable solution without issue. This highlights a key advantage of our method: its inherent stability makes it a reliable tool even for underdetermined problems, where many other regression techniques might fail or require regularization.

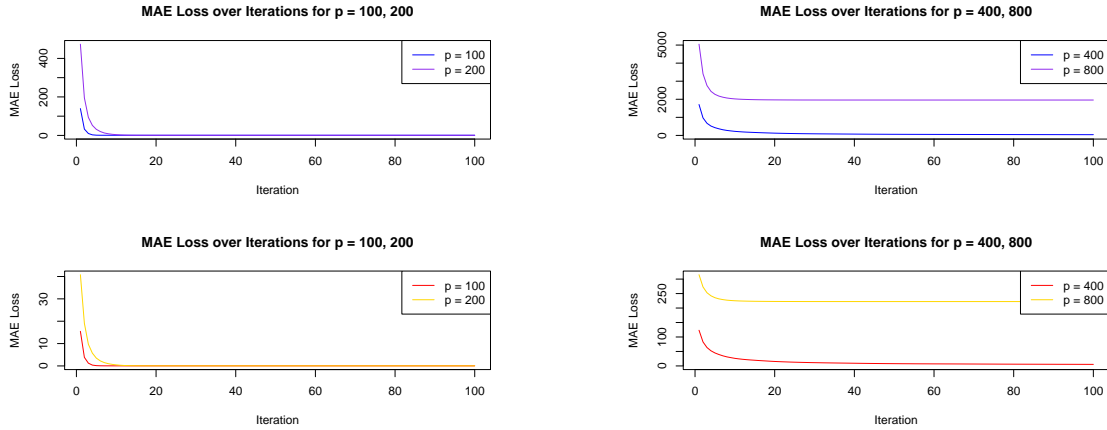


Figure 3: MAE on predictions over iterations for high-dimensional settings where  $p < n$ . Convergence is rapid in all cases, typically stabilizing in under 20 iterations.

### 4.3 Boston Housing dataset

To validate our proposed coordinate descent algorithm for LAD regression, we conducted experiments on the well-known Boston Housing dataset. This dataset consists of 506 observations and 13 covariates, with the response variable `medv` (median house value). Due to the presence of outliers and heteroskedasticity in this dataset, it serves as a standard benchmark for testing robust regression methods. We compared the performance of our coordinate descent implementation (denoted LAD-CD) with the classical `quantreg` package

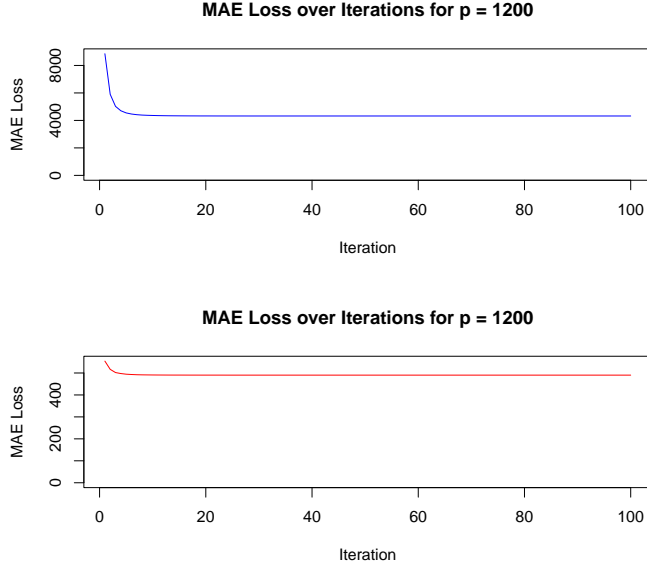


Figure 4: Algorithm performance in the under-determined case where  $p = 1200 > n = 1000$ . The algorithm remains stable and converges just as quickly, demonstrating its robustness.

(Koenker and d’Orey, 1987; Koenker and Bassett, 1978), which solves LAD regression using the Barrodale–Roberts linear programming method. Our primary comparison metrics were prediction accuracy (mean absolute error, MAE) and runtime efficiency.

Figure 5 [Left] shows predicted house values from both methods against the true values. The two models give very similar fitted values, demonstrating that LAD–CD converges to solutions close to those obtained from the well-established linear programming approach. Figure ?? [Right] tracks the MAE across iterations of the LAD–CD algorithm. The loss decreases monotonically, converging to a stable level close to the MAE achieved by `quantreg`. This confirms the theoretical convergence guarantees of our method.

Overall, while the LP-based solver is faster on small datasets, our LAD–CD implementation provides a simple and interpretable alternative that achieves comparable robustness and accuracy without relying on specialized optimization libraries.

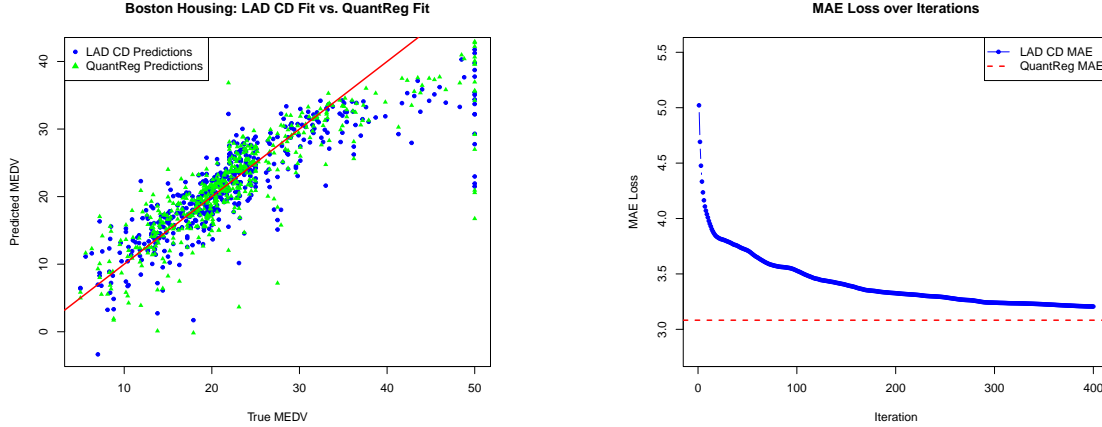


Figure 5: [Left] Predicted vs. true median house values for the Boston Housing dataset. Both LAD-CD (blue) and `quantreg` (green) produce similar fits, illustrating the consistency of our method. [Right] MAE loss over iterations for LAD-CD compared to `quantreg`. The coordinate descent method converges monotonically toward the optimal loss achieved by the LP-based solver.

Table 1: Comparing (LAD-CD) and (QuantReg) on the Boston Housing dataset

|                              | MAE (loss) | Runtime (s) | $\sum_j  \hat{\beta}_j^{CD} - \hat{\beta}_j^{rq} $ |
|------------------------------|------------|-------------|--|
| LAD-CD                       | 3.53       | 0.115       | 30.24  |
| QuantReg (Barrodale-Roberts) | 3.08       | 0.002       | —  |

#### 4.4 Air Quality: Comparison of LAD Coordinate Descent and Quantile Regression

We next evaluated our algorithm on the `airquality` dataset, which records daily air quality measurements in New York, including ozone concentration, solar radiation, wind, and temperature. After removing rows with missing values, we modeled ozone concentration as the response variable using all available predictors.

As in the Boston Housing case, we compared our coordinate descent implementation (LAD-CD) against the `quantreg` package, which employs the Barrodale-Roberts linear programming method. The goal of this experiment was to test robustness on a dataset with more variability and heavy-tailed residuals, where LAD regression is particularly well suited.

Figure 6 [Left] compares predicted ozone levels from both methods against the observed

values. Despite some scatter caused by extreme observations, both methods yield broadly consistent predictions. Figure 6 [Right] shows the monotone decrease of the LAD–CD loss over iterations, converging to a level close to the quantile regression baseline.

Numerical results are summarized in Table 2. Similar to the Boston dataset, LAD–CD achieves comparable predictive accuracy, though it requires more iterations and is slower than the optimized LP solver. The coefficient distance between the two methods is larger in this case, reflecting the greater variability and influence of outliers, but the overall predictive behavior remains consistent.

Table 2: Comparison of LAD Coordinate Descent (LAD–CD) and Quantile Regression (QuantReg) on the Air Quality dataset.

|                              | MAE (loss) | Runtime (s) | $\sum_j  \hat{\beta}_j^{CD} - \hat{\beta}_j^{rq} $ |
|------------------------------|------------|-------------|--|
| LAD–CD                       | 15.49      | 0.079       | 78.29  |
| QuantReg (Barrodale–Roberts) | 14.34      | 0.001       | —  |

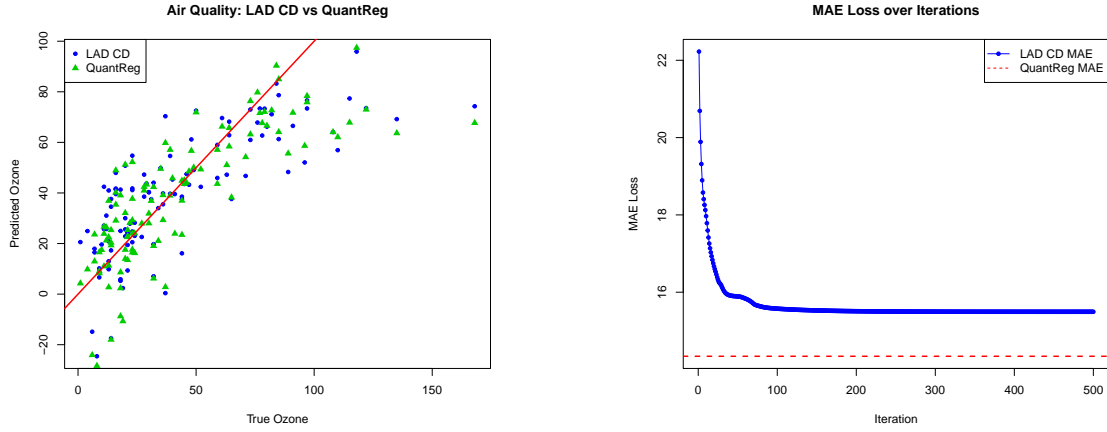


Figure 6: [Left] Predicted vs. true ozone levels for the Air Quality dataset. Both LAD–CD (blue) and QuantReg (green) yield broadly similar predictions, though the scatter reflects the heavy-tailed nature of the data. [Right] MAE loss over iterations for LAD–CD compared to QuantReg. The coordinate descent method converges monotonically, reaching a level close to the QuantReg baseline.

## 4.5 Concrete Compressive Strength: Comparison of LAD Coordinate Descent and Quantile Regression

As a third real-world example, we considered the `Concrete Compressive Strength` dataset from the UCI repository. This dataset contains 1,030 observations of concrete mixtures, with eight covariates describing material composition (e.g., cement, water, slag, fly ash, coarse and fine aggregates, superplasticizer, and age). The response variable is compressive strength, measured in MPa. This dataset is known to exhibit heteroskedasticity and outliers, making it a natural candidate for robust regression methods such as LAD.

We again compared our coordinate descent algorithm (LAD-CD) to the `quantreg` package’s Barrodale–Roberts solver. Both methods yielded very similar predictive accuracy, with MAEs of 8.14 and 8.05, respectively. Runtime results follow the same trend observed in previous experiments: the optimized LP-based solver is faster on this moderate-sized dataset, while LAD-CD remains efficient and simple to implement. The coefficient distance between the two fits is larger here, reflecting the greater variability of the predictors, but the predictive behavior remains consistent.

Table 3 summarizes the numerical results, while Figure 7 shows predicted versus observed strengths and the convergence of the LAD-CD loss.

Table 3: Comparison of LAD Coordinate Descent (LAD-CD) and Quantile Regression (QuantReg) on the Concrete Compressive Strength dataset.

|                              | MAE (loss) | Runtime (s) | $\sum_j  \hat{\beta}_j^{CD} - \hat{\beta}_j^{rq} $ |
|------------------------------|------------|-------------|--|
| LAD-CD                       | 8.14       | 0.258       | 98.91  |
| QuantReg (Barrodale–Roberts) | 8.05       | 0.003       | –  |

## 5 Conclusion

In this paper, we introduced a simple coordinate descent algorithm for solving the Least Absolute Deviations (LAD) regression problem. Motivated by the sensitivity of traditional OLS regression to outliers, our method provides a robust alternative that is easy to implement without the overhead of linear programming solvers. By iteratively solving one-dimensional subproblems with closed-form solutions based on medians and weighted medians, our approach is guaranteed to converge to the global optimum.

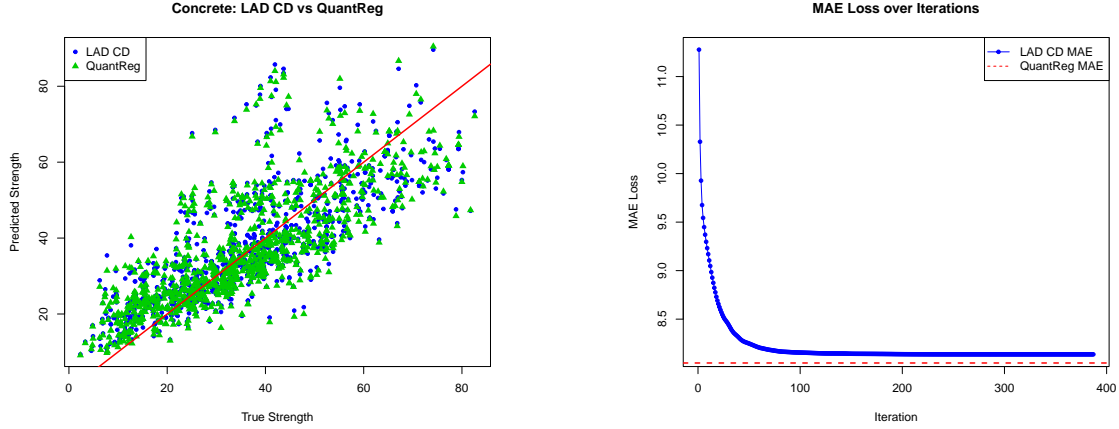


Figure 7: [Left] Predicted vs. true compressive strengths for the Concrete dataset. Both LAD-CD (blue) and QuantReg (green) yield similar predictive accuracy. [Right] MAE loss over iterations for LAD-CD compared to QuantReg. The coordinate descent method converges monotonically toward the quantile regression baseline.

Theoretical analysis established the convexity of the LAD objective and the validity of coordinate-wise updates, while experiments on synthetic and real-world data confirmed the practical utility of the method. Across three representative datasets—Boston Housing, Air Quality, and Concrete Compressive Strength—LAD-CD achieved predictive performance comparable to the classical `quantreg` package, with nearly identical MAE values. Although the linear programming solver was faster in absolute runtime, our method consistently converged monotonically and produced stable fits, even in noisy or high-dimensional settings.

Overall, the proposed algorithm demonstrates that LAD regression can be made accessible and transparent with only a few lines of code. This makes it an attractive option for teaching, for quick prototyping, and for applied problems where robustness is critical but specialized optimization software may not be available. Future extensions could explore penalized LAD regression, parallel implementations to further improve scalability, and applications to larger and more complex datasets.

## References

- Barrodale, I. and Roberts, F. D. K. (1973). An improved algorithm for discrete  $l_1$  linear approximation. *SIAM Journal on Numerical Analysis*, 10(5):839–848.
- Barrodale, I. and Roberts, F. D. K. (1974). Solution of an overdetermined system of equations in the  $l_1$  norm. *Communications of the ACM*, 17(6):319–320. Algorithm 478.
- Charnes, A., Cooper, W. W., and Ferguson, R. O. (1955). Optimal estimation of executive compensation by linear programming. *Management Science*, 1(2):138–151.
- Koenker, R. and Bassett, Gilbert, J. (1978). Regression quantiles. *Econometrica: Journal of the Econometric Society*, pages 33–50.
- Koenker, R. W. and d’Orey, V. (1987). Computing regression quantiles. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 36(3):383–393.
- Mkhadri, A., Ouhourane, M., and Oualkacha, K. (2017). A coordinate descent algorithm for computing penalized smooth quantile regression. *Statistics and Computing*, 27:865–883.
- Peng, B. and Wang, L. (2015). An iterative coordinate descent algorithm for high-dimensional nonconvex penalized quantile regression. *Journal of Computational and Graphical Statistics*, 24(3):676–694.
- Wagner, H. M. (1959). Linear programming techniques for regression analysis. *Journal of the American Statistical Association*, 54(285):206–212.