# Advance Regression Term Project

Gerard Palomo & Juan Pablo Uphoff

**Abstract**
Ordinary least squares (OLS) regression fits the mean and is sensitive to extreme observations, whereas quantile regression (e.g. median regression) fits conditional quantiles and tends to be more robust to outliers. We conduct a simulation study comparing OLS and median (50% quantile) regression under increasing levels of outlier contamination in the response variable. Synthetic data are generated from a simple linear model and a fraction of responses are replaced by large aberrant values. We fit both OLS and quantile regression and compare the resulting coefficient estimates. As expected, OLS estimates become severely biased in the presence of outliers, while the median regression estimates remain closer to the true values. These findings confirm that quantile regression can better resist outliers than OLS

*Corresponding author:

## Contents

## 1. Introduction.

Classical linear regression (ordinary least squares, OLS) fits a line by minimizing the sum of squared residuals, which targets the conditional mean of the response variable. However, because OLS puts heavy weight on large deviations, its estimates can be strongly affected by even a few outliers. In contrast, quantile regression extends median regression to any conditional quantile (e.g. the 50% quantile) and minimizes an asymmetrically weighted absolute loss. A key advantage of quantile regression is its robustness: the influence function of the median is bounded, so regression quantiles inherit a high breakdown point. In practical terms, estimates like the conditional median are not pulled as far by extreme values. Prior studies confirm these properties: for example, Herawati (2020) showed in simulations that median regression provided smaller mean-squared error than OLS when outliers are present. In this article, we compare OLS and quantile regression for a linear model under controlled outlier contamination.

## 2. Linear Quantile Regression: Theory and Methods

### 2.1 Definition and Estimation

For a random variable $Y$, the $\tau$-th quantile is the value $q_\tau$ such that $P(Y \leq q_\tau) = \tau$. In a regression setting, **quantile regression (QR)** estimates:

$$Q_Y(\tau \mid X = x) = x^\top \beta(\tau),$$

where $\beta(\tau)$ is a vector of coefficients specific to quantile level $\tau$. For example, $\beta_1(0.5)$ represents the effect of $X_1$ on the median of $Y$.

Koenker and Bassett (1978) proposed estimating $\beta(\tau)$ by minimizing the **check loss**:

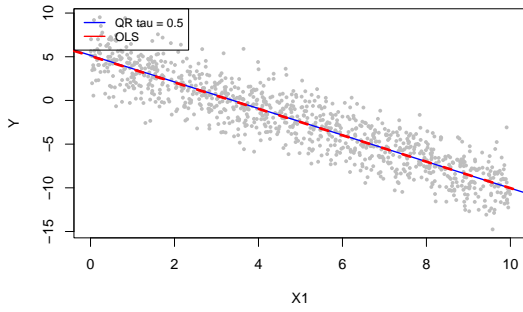$$\hat{\beta}(\tau) = \arg\min_\beta \sum_{i=1}^n \rho_\tau(y_i - x_i^\top \beta),$$

with $\rho_\tau(u) = u(\tau - \mathbb{I}\{u < 0\})$. For $\tau = 0.5$, this reduces to **least absolute deviations (LAD)** regression.

Each $\tau$ is estimated independently using linear programming, and the family $\{\beta(\tau)\}$ forms a **quantile process** describing the full conditional distribution of $Y$.

### 2.2 Inference

Under regularity conditions, $\hat{\beta}(\tau)$ is asymptotically normal:

$$\sqrt{n}(\hat{\beta}(\tau) - \beta(\tau)) \xrightarrow{d} N(0, \Sigma(\tau)),$$

**Figure 1.** OLS vs Multiple Quantile Regressions (Clean Data)



**Figure 2.** OLS vs Quantile Regressions (Outliers Data)

**Table 1.** Comparison of OLS and Quantile Regression Coefficients (Clean Data) - Manual

| Coefficient | OLS | QR 0.50 |
|---|---|---|
| (Intercept) | 5.082 | 5.16 |
| X1 | -1.512 | -1.52 |

**Table 2.** Comparison of OLS and Quantile Regression Coefficients (Clean Data) - Manual

| Coefficient | OLS | QR 0.50 |
|---|---|---|
| (Intercept) | 0.390 | 4.723 |
| X1 | -0.065 | -1.391 |



**Figure 3.** Slope vs Quantile for heteroskedastic case

- **Homoscedastic Gaussian**: $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, $\sigma = 2$. This setting satisfies OLS assumptions and serves as a benchmark.
- **Contaminated Gaussian (outliers)**: $\varepsilon_i = Z_i + O_i$ with $Z_i \sim \mathcal{N}(0, 2^2)$ and $O_i = 50$ with probability $p = 0.05$. This mimics a mixture contamination model, testing robustness. Outliers are introduced conditionally (to the largest $X_1$) to induce high-leverage effects.
- **Heteroskedastic**: $\varepsilon_i \sim \mathcal{N}(0, \sigma^2(X_i))$ with $\sigma(X_i) = 1 + 0.2X_1$. This setting creates increasing spread in $Y$ conditional on $X_1$, violating OLS homoscedasticity.
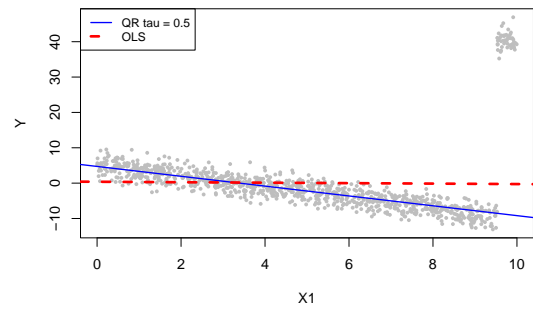
### 3.3 Evaluation Metrics

*****PENDING***** Along with some graphical support, these evaluations provide insight into estimator consistency, robustness, and efficiency under both ideal and adversarial conditions.

The theoretical basis for QR (Section 2) suggests that: - OLS is optimal under homoscedastic Gaussian errors (Gauss-Markov assumptions), - QR provides robust, distribution-sensitive estimation that remains valid when these assumptions are violated.
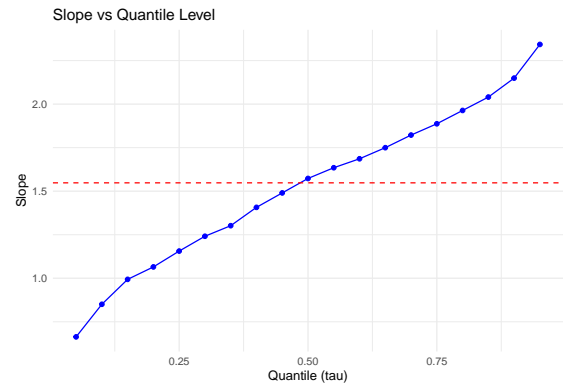
This simulation seeks to verify those expectations empirically.

### 3.4 Simulation

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim

proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.
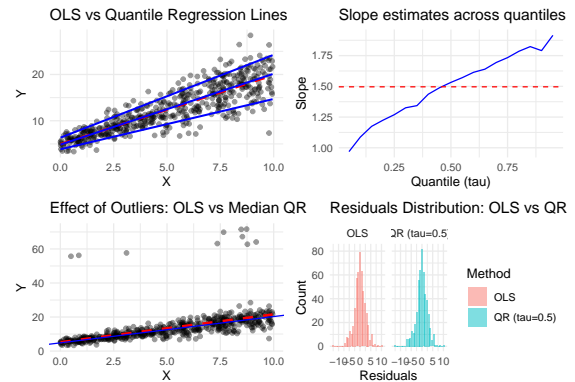


**Figure 4.** Various plots

## 4. Discussion & Conclusions.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.