

Lecture 10

Robust and Quantile Regression

1

Outliers

- Many definitions: Atypical observations, extreme values, conditional unusual values, observations outside the expected relation, etc.
- In general, we call an *outlier* an observation that is numerically different from the data. But, is this observation a “mistake,” say a result of measurement error, or part of the (heavy-tailed) distribution?
- In the case of normally distributed data, roughly 1 in 370 data points will deviate from the mean by $3xSD$. Suppose $T=1000$. Then, 9 data points deviating from the mean by more than $3xSD$ indicates outliers. But, which of the 9 observations can be classified as an outliers?
- Problem with outliers: They can affect estimates. For example, with small data sets, one big outlier can seriously affect OLS estimates.

Outliers

- Several identifications methods:

- *Eyeball*: Look at the observations away from a scatter plot.

- *Standardized residual*: Check for errors that are two or more standard deviations away from the expected value.

- *Leverage statistics*: It measures the difference of an independent data point from its mean. High leverage observations can be potential outliers. Leverage is measured by the diagonal values of the **P** matrix:

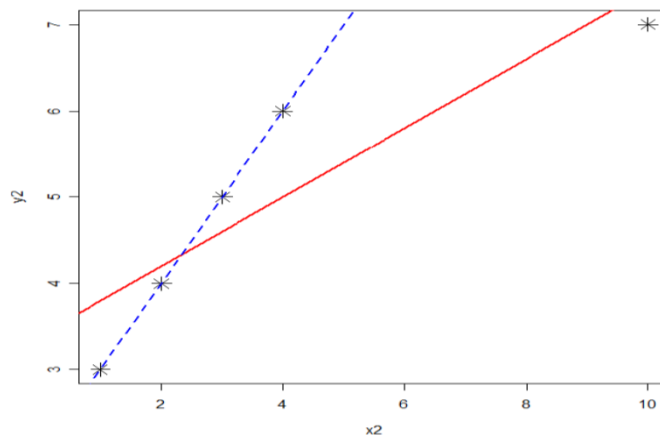
$$h_i = 1/T + (x_i - \bar{x})/[(T-1)s_x^2].$$

But, an observation can have high leverage, but no *influence*.

- *Influence statistics: Dif beta*. It measures how much an observation influences a parameter estimate, say b_j . Dif beta is calculated by removing an observation, say i , recalculating b_j , say $b_{j(-i)}$, taking the difference in betas and standardizing it. Then,

$$Dif\ beta_{j(-i)} = [b_j - b_{j(-i)}]/SE[b_j].$$

Outliers



- Deleting the observation in the upper right corner has a clear effect on the regression line. This observation has *leverage* and *influence*.

Outliers

- A related popular influence statistic is *Distance D (as in Cook's D)*. It measures the effect of deleting an observation on the fitted values, say \hat{y}_j

$$D_j = \sum_i [\hat{y}_j - \hat{y}_j(-i)] / [K \text{MSE}]$$

where K is the number of parameters in the model and MSE is mean square error of the regression model.

- The influence statistics are usually compare to some ad-hoc cut-off values used for identifying highly influential points, say $D_i > 4/T$.
- The analysis can also be carried out for groups of observations. In this case, we would be looking for blocks of highly influential observations.

Outliers: Summary of Rules of Thumb

- General rules of thumb used to identify outliers:

Measure	Value
abs(stand resid)	> 2
abs(Dif Beta)	$> 2/\text{sqrt}(T)$
Cook's D	$> 4/T$
leverage	$> (2k+2)/T$

Outliers: Regression – SAS - Application

```
proc reg data = ab;
model S51-RF = MKT-RF SMB HML / white vif collinoint;
output out=capmres(keep=year S51-RF MKT-RF SMB HML
r sr lev cd dffit)
r=res student=sr h=lev cookd=cd;
run;
```

The REG Procedure
Model: MODEL1
Dependent Variable: S51-RF

Root MSE	3.44666	R-Square	0.8857
Dependent Mean	0.99706	Adj R-Sq	0.8853

Parameter Estimates

Variable	Label	DF	Parameter Estimate	Standard Error	t Value	White SE
Intercept	Intercept	1	-0.33505	0.10816	-3.10	0.09035
xm		1	1.03766	0.02128	48.76	0.03982
SMB	SMB	1	1.51900	0.03441	44.15	0.09993
HML	HML	1	0.74036	0.03095	23.92	0.08977

Outliers: Distribution – SAS - Application

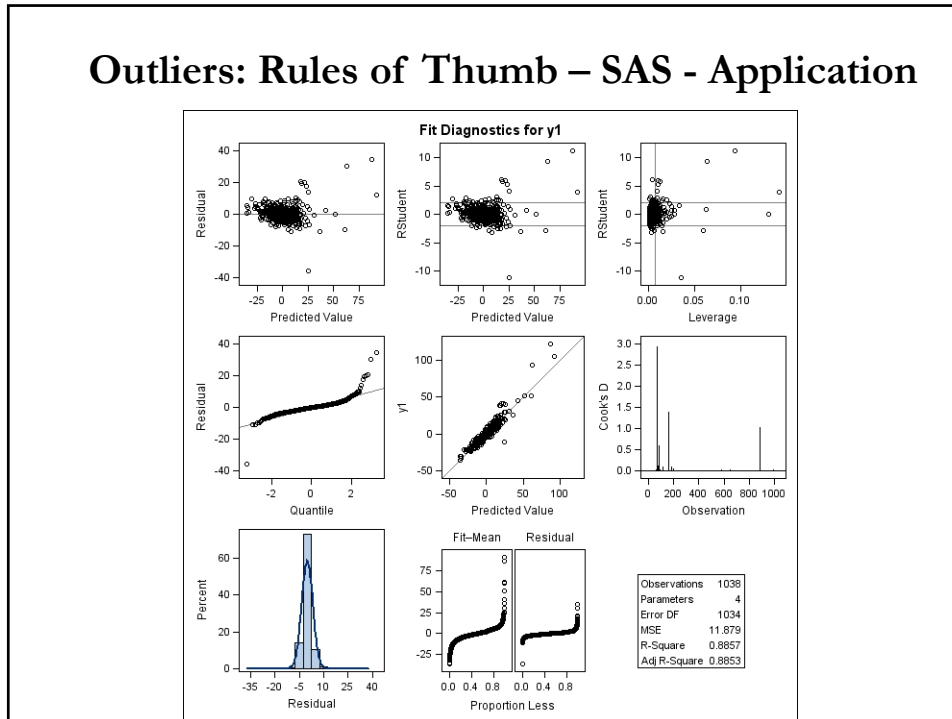
The UNIVARIATE Procedure
Variable: r (Studentized Residual without Current Obs)

Basic Statistical Measures

Location		Variability	
Mean	0.00097	Std Deviation	1.00494
Median	-0.09766	Variance	1.00990
Mode	.	Range	8.34790
		Interquartile Range	1.16010

Quantile	Estimate
100% Max	4.7696676
95%	1.7128005
90%	1.2168926
75% Q3	0.5638215
50% Median	-0.0976612
25% Q1	-0.5962799
10%	-1.1582571
5%	-1.4562294
0% Min	-3.5782300

Outliers: Rules of Thumb – SAS - Application



Outliers

- What to do?
 - Use a non-linear formulation or apply a transformation (log, square root, etc.) to the data.
 - Remove suspected observations. (Sometimes, there are theoretical reasons to remove suspect observations. Typical procedure in finance, remove public utilities or financial firms from the analysis.)
 - Winsorization of the data.
 - Use dummy variables.
 - Use LAD (quantile) regressions, which are less sensitive to outliers.
 - Weight observations by size of residuals or variance (robust estimation).

- General rule: Present results with or without outliers.

Robust Estimation

- Following Huber (1981), we will interpret *robustness* as insensitivity to small deviations from the assumptions the model imposes on the data.
- In particular, we are interested in *distributional robustness*, and the impact of skewed distributions and/or *outliers* on regression estimates.
 - In this context, *robust* refers to the shape of a distribution –i.e., when the actual distribution differs from the theoretically assumed distribution.
 - Although conceptually distinct, distributional robustness and outlier resistance are, for practical purposes, synonymous
 - Robust can also be used to describe standard errors that are adjusted for non-constant error variance. But, we have already covered this topic.

13

Robust Estimation – Mean vs Median

- Intuition: Under normality, OLS has optimal properties. But, under non-normality, nonlinear estimators may be better than LS estimators.

Example: *i.i.d.* case

Let $\{y_i\}_1^T \sim iid F\left(\frac{y-\mu}{\sigma}\right)$ where $F(0) = 0.5$

where F is a symmetric distribution with scale parameter σ .

- Let the order statistics be $y_{(1)} \leq \dots \leq y_{(T)}$
- Sample median: $\tilde{\mu} = y_{\left(\frac{T+1}{2}\right)}$
- Laplace showed that

$$\sqrt{T}(\tilde{\mu} - \mu) \rightarrow N\left(0, \frac{1}{4f(\mu=0)^2}\right)$$

14

Robust Estimation – Mean vs Median

- Using this result, one can show:

	$T \text{ var}(\text{mean})$	$T \text{ var}(\tilde{\mu} = \text{median})$
Normal	1	1.57
Laplace	2	1
Average	1.5	1.28

- Intuitively, this occurs because Laplace is fat-tailed, and the median is much less sensitive to the information in the tails than the mean.
- The mean gives $1/T$ weight to all observations (close to the mean or in the tails). A large observation can seriously affect (*influence*) the mean, but not the median.

15

Robust Estimation – Mean vs Median

- Remark: The sample mean is the MLE under the Normal distribution; while the sample median is the MLE under the Laplace distribution.

- If we do not know which distribution is more likely, following Huber, we say the median is robust (“better”). But, if the data is normal, the median is not efficient (57% less efficient than mean).
- There are many types of robust estimators. Although they work in different ways, they all give less weight to observations that would otherwise influence the estimator.
- Ideally, we would like to design a weighting scheme that delivers a robust estimator with good properties (efficiency) under normality. ¹⁶

Robust Estimation – Mean vs Median

Examples: Robust estimators for central location parameter.

- The sample median, $\tilde{\mu}$.
- Trimmed-Mean, the mean of the sample after fraction α of the largest and smallest observations have been removed.
- The “Winsorized Mean:”

$$\hat{\mu}^w = \frac{1}{T} \left((g+1)y_{(g+1)} + y_{(g+2)} + \dots + y_{(T-g-1)} + (g+1)y_{(T-g)} \right)$$

which is similar to the trimmed-mean, but instead of throwing out the extremes, we “accumulate” them at the truncation point.

- Q: All robust, which one is better? Trade-off: robustness-efficiency.
- The concept of robust estimation can be easily extended to the problem of estimating parameters in the regression framework.

17

Robust Regression

- There are many types of robust regression models. Although they work in different ways, they all give less weight to observations that would otherwise influence the regression line.

- Early methods:

– **Least Absolute Deviation/Values** (LAD/LAV) regression or least absolute deviation regression –i.e., minimizes $|e|$ instead of e^2 .

- Modern methods:

- ***M-Estimation***

- Huber estimates, Bi-square estimators

- ***Bounded Influence Regression***

- Least Median of Squares, Least-Trimmed Squares

18

Review: M-Estimation

- An extremum estimator is one obtained as the optimizer of a criterion function, $q(\mathbf{z}, \mathbf{b})$.

Examples:

$$\text{OLS: } \mathbf{b} = \arg \max (-\mathbf{e}'\mathbf{e}/T)$$

$$\text{MLE: } \mathbf{b}_{\text{MLE}} = \arg \max \ln L = \sum_{i=1, \dots, T} \ln f(y_i, \mathbf{x}_i, \mathbf{b})$$

- M-estimators: The objective function is a sample average or a sum.
 - "M" stands for a maximum or minimum estimators --Huber (1967). It can be viewed a generalization of MLE.
- We want to obtain: $\mathbf{b}_M = \operatorname{argmin} \sum_i q(\mathbf{z}_i, \mathbf{b})$ (or divided by T).
 - If $q(y_i - \mathbf{x}_i' \mathbf{b}_M)$, $q(\cdot)$ measures the contribution of each residual to the objective function.

19

Review: M-Estimation

- We want to obtain: $\mathbf{b}_M = \operatorname{argmin} \sum_i q(y_i - \mathbf{x}_i' \mathbf{b}_M)$, (or divided by T)
- In general, we solve the f.o.c. Let $\psi = \partial q(\cdot) / \partial \mathbf{b}'$. Then,

$$\sum_i \psi(y_i - \mathbf{x}_i' \mathbf{b}_M) \mathbf{x}_i' = \mathbf{0} \quad (\text{K equations})$$

- We replace $\psi(\cdot)$ with the weight function, $w_i = \psi(e_i) / e_i$

$$\sum_i w_i (y_i - \mathbf{x}_i' \mathbf{b}) \mathbf{x}_i = \sum_i w_i e_i \mathbf{x}_i = \mathbf{0}$$

These f.o.c.'s are equivalent to a weighted LS problem, which minimizes $\sum_i w_i e_i^2$.

- Q: Which $q(\cdot)$, or equivalently, w_i should we use to produce a *robust* estimator?

20

M-Estimation: Asymptotic Normality

- Summary
 - $\mathbf{b}_M \xrightarrow{p} \mathbf{b}_0$
 - $\mathbf{b}_M \xrightarrow{a} N(\mathbf{b}_0, \text{Var}[\mathbf{b}_0])$
 - $\text{Var}[\mathbf{b}_M] = (1/T) \mathbf{H}_0^{-1} \mathbf{V}_0 \mathbf{H}_0^{-1}$
 - If the model is correctly specified: $-\mathbf{H} = \mathbf{V}$.
 - Then, $\text{Var}[\mathbf{b}] = \mathbf{V}_0$
 - \mathbf{H} and \mathbf{V} are evaluated at \mathbf{b}_0 :
 - $\mathbf{H} = \sum_i [\partial^2 q(\mathbf{z}_i, \mathbf{b}) / \partial \mathbf{b} \partial \mathbf{b}']$
 - $\mathbf{V} = \sum_i [\partial q(\mathbf{z}_i, \mathbf{b}) / \partial \mathbf{b}] [\partial q(\mathbf{z}_i, \mathbf{b}) / \partial \mathbf{b}]'$

21

M-Estimators in the Regression Context

- Many $q(\mathbf{z}, \beta)$ can be structured to deliver a *robust* estimator?
- For example, we can define the family of L_p -estimators:
 - $q(\mathbf{z}; \beta) = (1/p) |\mathbf{x} - \beta|^p$ for $1 \leq p \leq 2$
 - $\mathbf{s}(\mathbf{z}; \beta) = |\mathbf{x} - \beta|^{p-1}$ $\mathbf{x} - \beta < 0$
 - $= -|\mathbf{x} - \beta|^{p-1}$ $\mathbf{x} - \beta > 0$
- Special cases:
 - $p = 2$: We get the sample mean (LS estimator for β).
 - $\mathbf{s}(\mathbf{z}; \beta) = \sum_i (x_i - \mathbf{b}_M) = 0 \quad \Rightarrow \quad \mathbf{b}_M = \sum_i x_i / T$

$p = 1$: We get the sample median as the estimator with the least absolute deviation (LAD) for the median β . (No unique solution if T is even.). Numerical (linear programming) solution needed.

22

M-Estimators in the Regression Context: Example

```

-----
Least absolute deviations estimator.....
Residuals  Sum of squares      =    1537.58603
            Standard error of e =     6.82594
Fit        R-squared        =     .98284
            Adjusted R-squared =     .98180
Sum of absolute deviations      =    189.3973484
-----
Variable| Coefficient   Standard Error  b/St.Er.  P[|Z|>z]  Mean of X
-----|-----
Covariance matrix based on 50 replications.
Constant| -84.0258***    16.08614      -5.223    .0000
Y        |  .03784***     .00271        13.952    .0000    9232.86
PG       | -17.0990***    4.37160       -3.911    .0001    2.31661
-----
Ordinary least squares regression .....
Residuals  Sum of squares      =    1472.79834
            Standard error of e =     6.68059  Standard errors are based on
Fit        R-squared        =     .98356  50 bootstrap replications
            Adjusted R-squared =     .98256
-----
Variable| Coefficient   Standard Error  t-ratio  P[|T|>t]  Mean of X
-----|-----
Constant| -79.7535***    8.67255       -9.196    .0000
Y        |  .03692***     .00132        28.022    .0000    9232.86
PG       | -15.1224***    1.88034       -8.042    .0000    2.31661
-----

```

23

Breakdown Point: Intuition

- There are several measures of robustness of an estimator, attempting to quantify the change. One of the most commonly used is the *breakdown point*.
- Let \mathbf{X} be a random sample and $\mathbf{T}(\mathbf{X})$ be an estimator. Informally, the breakdown point of the estimator is the proportion m/T of observations, which can be replaced by *bad observations* (outliers) without forcing $\mathbf{T}(\mathbf{X})$ to leave a bounded set –i.e., become infinity.

Example: The sample mean has a breakdown point equal to 0 (one observation can drive the sample mean, regardless of the other $T-1$ values). The median has a breakdown point $1/2$ (it can tolerate 50% bad values) and $\alpha\%$ -trimmed mean has a breakdown point $\alpha\%$.

24

Breakdown Point: Definition

- Assume a sample, \mathbf{Z} , with T observations, and let \mathbf{T} be a regression estimator. That is, we apply \mathbf{T} to \mathbf{Z} we get the regression coefficients:

$$\mathbf{T}(\mathbf{Z}) = \mathbf{b}$$

- Imagine all possible “corrupted” samples \mathbf{Z}^0 that replace any subset of observations, m , in the dataset with arbitrary values -*i.e.*, influential cases.

- The maximum bias that could arise from these substitutions is:

$$\text{bias}(m; \mathbf{T}, \mathbf{Z}) = \sup_{\mathbf{Z}^0} \|\mathbf{T}(\mathbf{Z}^0) - \mathbf{T}(\mathbf{Z})\|$$

- If the $\text{bias}(m; \mathbf{T}, \mathbf{Z})$ is infinite, the m outliers have an arbitrarily large effect on \mathbf{T} . In other words, the estimator *breaks down*.

25

Breakdown Point: Definition

- Then, the breakdown point for an estimator \mathbf{T} for a finite sample \mathbf{Z} is:

$$\varepsilon_n^*(\mathbf{T}, \mathbf{Z}) = \min \left\{ \frac{m}{n}; \text{bias}(m; \mathbf{T}, \mathbf{Z}) \text{ is infinite} \right\}$$

- The breakdown point of an estimator is the smallest fraction of “*bad*” data (outliers or data grouped at the extreme of a tail) the estimator can tolerate without taking on values arbitrarily far from $\mathbf{T}(\mathbf{Z})$.

- For OLS regression one unusual case is enough to influence the coefficient estimates. Its breakdown point is then

$$\varepsilon_n^*(\mathbf{T}, \mathbf{Z}) = 1/T$$

- As T gets larger, $1/T$ tends towards 0, meaning that the breakdown point for OLS is 0%.

26

Robust Regression: Methods

- Robust regression methods attempt to limit the impact of unusual cases on the regression estimates
 - **Least Absolute Values (LAV/LAD) regression** is robust to outliers (unusual Y values given X), but typically fares even worse than OLS for cases with high leverage.
 - If a leverage point is very far away, the LAD line will pass through it. In other words, its breakdown point is also $1/T$.
 - **M-Estimators** are also robust to outliers. More efficient than LAD estimators. They can have trouble handling cases with high leverage, meaning that the breakdown point is also $1/T$.
 - **Bounded influence methods** have a much higher breakdown point (as high as 50%) because they effectively remove a large proportion of the cases. These methods can have trouble with small samples.

27

Estimating the Center of a Distribution

- In order to explain how robust regression works, we start with the simple case of robust estimation of the centre of a distribution. Consider independent observations and the simple model:

$$Y_i = \mu + \epsilon_i$$

- If the underlying distribution is normal, the sample mean is the MLE.
- The mean minimizes the LS objective function:

$$q_{LS} = \mathbf{e}'\mathbf{e} = \sum_i e_i^2$$

- The derivative of the objective function with respect to e_i gives the influence function which determines the influence of observations: $\psi_{LS,i}(e) = 2 e_i$. That is, influence is proportional to the residual e_i .

28

Estimating the Center of a Distribution

- As an alternative to the mean, we consider the median as an estimator of μ . The median minimizes the LAD objective function:

$$Q_{\text{LAD}} = 1/T \sum_i |e_i|$$

- Taking the derivative of the objective function gives the shape of the influence function:

$$\begin{aligned} \psi_{\text{LAD},i}(e) &= 1 && \text{for } e_i > 0. \\ &= 0 && \text{for } e_i = 0. \\ &= -1 && \text{for } e_i < 0. \end{aligned}$$

- Note that influence of e_i is bounded. The fact that the median is more resistant than the mean to outliers is a favorable characteristic.

29

Influence Function for Mean and Median

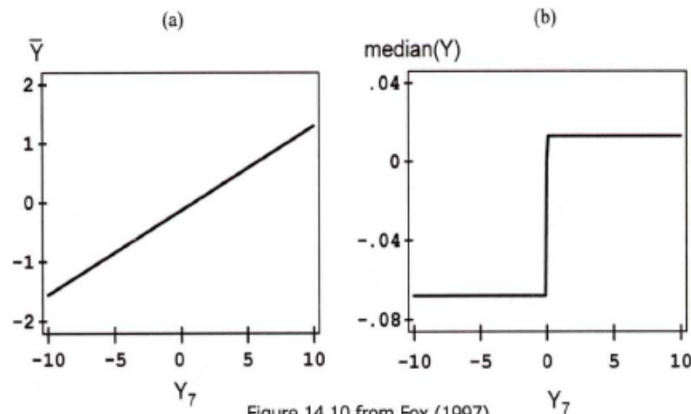


Figure 14.10 from Fox (1997)

30

M-Estimation: Huber Estimates

- But, the median is far less efficient, however. If $Y \sim N(\mu, \sigma^2)$,

$$\text{Var}[\bar{\mu}] = \sigma^2/T$$

$$\text{Var}[\tilde{\mu}] = \pi\sigma^2/2T$$

=> The $\text{Var}[\tilde{\mu}]$ is $\pi/2$ (≈ 1.57) times as large as $\text{Var}[\text{mean}]$.

- A good compromise between the efficiency of LS and the robustness of LAD is the Huber (1964) objective function:

$$\begin{aligned} \rho_{H,i}(e_i) &= \frac{1}{2} e_i^2 && \text{for } |e_i| \leq k. && (k = \text{tuning constant}) \\ &= k |e_i| - \frac{1}{2} k^2 && \text{for } |e_i| > k. \end{aligned}$$

with an influence function:

$$\begin{aligned} \psi_{H,i}(e_i) &= k && \text{for } e_i > k. \\ &= e_i && \text{for } |e_i| \leq k. \\ &= -k && \text{for } e_i < -k. \end{aligned}$$

31

M-Estimation: Tuning constant, k

- k is called the *tuning constant*.

Note: For $k \rightarrow \infty$, the M-estimator turns into mean, for $k \rightarrow 0$, it becomes the median.

- Assuming the $\sigma=1$, setting $k=1.345$ produces 95% efficiency relative to the sample mean when the population is normal and gives substantial resistance to outliers when it is not.
- In general, k is expressed as a multiple of the *scale* of Y (the spread), S
=> $k=cS$.
 - We could use σ as a measure of scale, but it is more influenced by extreme observations than is the mean.
 - Instead, we use the *median absolute deviation*:

$$\text{MAD} = \text{median} |Y_i - \tilde{\mu}| = \text{median} |e_i|$$

32

M-Estimation: Tuning constant, k

- We use the *median absolute deviation*:

$$\text{MAD} = \text{median} |Y_i - \tilde{\mu}| = \text{median} |e_i|$$

- The median of Y serves as an initial estimate of $\tilde{\mu}$, thus allowing us to define $S = \text{MAD} / .6745$, which ensures that S estimates σ when the population is normal –i.e., for the standard normal $E[\text{MAD}] = 0.6745$
- Using $k = 1.345 S$ ($1.345 / .6745$ is about 2) produces 95% efficiency relative to the sample mean when the population is normal and gives substantial resistance to outliers when it is not.

Note: A smaller k gives more resistance to outliers.

33

M-Estimation: Bi-weight Estimates

- Tukey's *bi-weight (bisquare) estimates* behave somewhat differently than Huber weights, but are calculated in a similar manner
- The *biweight objective function* is especially resistant to observations on the extreme tails:

$$\begin{aligned} \rho_{\text{BW},i}(e_i) &= k^2/6 \{1 - [1 - (e_i/k)^2]^3\} && \text{for } |e_i| \leq k. \\ &= k^2/6 && \text{for } |e_i| > k. \end{aligned}$$

with an influence function:

$$\begin{aligned} \psi_{\text{BW},i}(e_i) &= \{e_i [1 - (e_i/k)^2]^2\} && \text{for } |e_i| \leq k. \\ &= 0 && \text{for } |e_i| > k. \end{aligned}$$

- For this function, $k = 4.685 S$ ($4.685 / .6745$ about 7 MADs) produces 95% efficiency when sampling from a normal population

34

M-Estimation and Regression

- Since regression is based on the mean, it is easy to extend the idea of M-estimation to regression. The linear model is:

$$y_i = \mathbf{x}_i' \mathbf{b} + \varepsilon_i$$

- The M-estimator then minimizes the objective function:

$$q = \sum_i q(y_i - \mathbf{x}_i' \mathbf{b})$$

with f.o.c.'s:

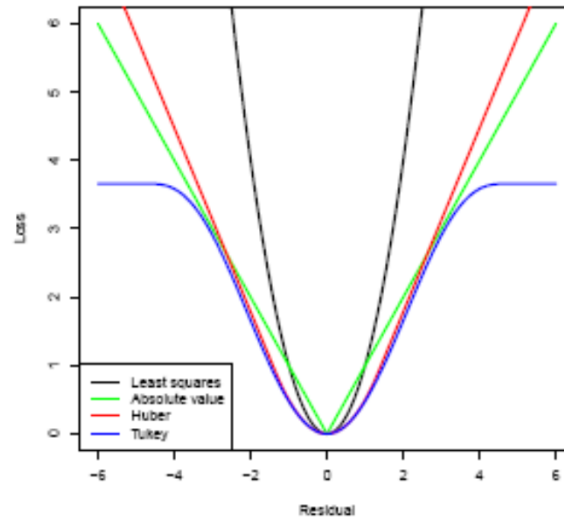
$$\sum_i \psi(y_i - \mathbf{x}_i' \mathbf{b}) \mathbf{x}_i' = \mathbf{0}$$

- We have a system of K equations. We replace $\psi(\cdot)$ with the weight function, $w(\varepsilon_i) = \psi(\cdot)/\varepsilon_i$: $\sum_i w_i (y_i - \mathbf{x}_i' \mathbf{b}) \mathbf{x}_i' = \mathbf{0}$

- The solution assigns a different weight to each case depending on the size of their residual; similar to a weighted least squares problem. ³⁵

M-Estimation and Regression: Loss Functions

- Different loss functions:



36

M-Estimation and Regression: Weights

- The weight function: $w(e_i) = \psi(\cdot)/e_i$;

Method	Objective Function	Weight Function
Least-Squares	$\rho_{LS}(e) = e^2$	$w_{LS}(e) = 1$
Huber	$\rho_H(e) = \begin{cases} \frac{1}{2}e^2 & \text{for } e \leq k \\ k e - \frac{1}{2}k^2 & \text{for } e > k \end{cases}$	$w_H(e) = \begin{cases} 1 & \text{for } e \leq k \\ k/ e & \text{for } e > k \end{cases}$
Bisquare	$\rho_B(e) = \begin{cases} \frac{k^2}{6} \left\{ 1 - \left[1 - \left(\frac{e}{k} \right)^2 \right]^3 \right\} & \text{for } e \leq k \\ k^2/6 & \text{for } e > k \end{cases}$	$w_B(e) = \begin{cases} \left[1 - \left(\frac{e}{k} \right)^2 \right]^2 & \text{for } e \leq k \\ 0 & \text{for } e > k \end{cases}$

37

Weight Functions for Various Estimators

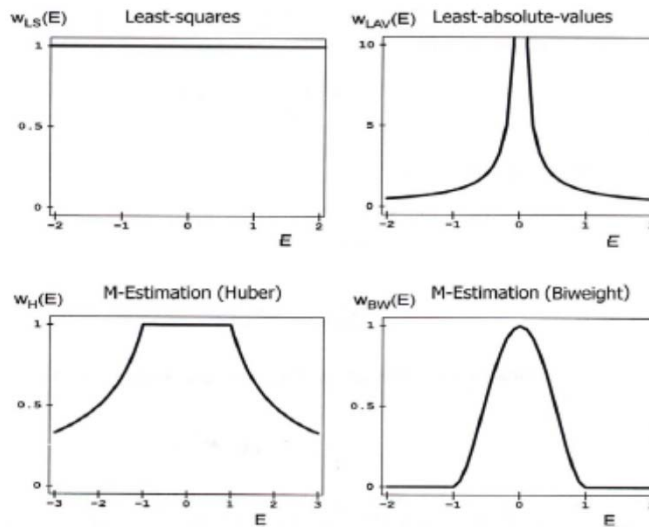


Figure 14.13 from Fox (1997)

38

M-Estimation and Regression: Algorithm

- The solution assigns a different weight to each case depending on the size of their residual, and thus minimizes the weighted sum of squares.

$$\sum_i w_i e_i^2 = 0$$

- The w_i weights depend on the residuals in the model. An iterative solution (using *Iterative Re-weighted Least Squares*, IRLS) is needed.

- The solution to this problem is weighted LS:

(1) Set initial \mathbf{b}^0 , say by using OLS. Get e_i^0 .

(2) Estimate the scale of the residuals S^0 and the weights w_i^0 .

(3) Estimate \mathbf{b}^j : j=1,2,...

$$\mathbf{b}^j = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}\mathbf{y} \quad \mathbf{W} = \text{diag}\{w_i^{j-1}\}$$

(4) With \mathbf{b}^j go back to (1). Repeat steps (1)-(3) until convergence.

39

M-Estimation and Regression

- Usual weight functions: Huber and Biweight (bisquare) weights.
- M-Estimators are statistically equally efficient as OLS if the distribution is normal, while at the same time are more robust with respect to influential cases.
- However, M-estimation can still be influenced by a single very extreme X-value—*i.e.*, like OLS, it still has a breakdown point of 0

40

Bounded Influence Regression: LTS

- M-estimation can still be influenced by a single very extreme X-value—*i.e.*, like OLS, it still has a breakdown point of 0
- *Least-trimmed-squares* (LTS) estimators –see Rousseeuw (1984)- can have a breakdown point up to 50% -*i.e.*, half the data can be influential in the OLS sense before the LTS estimator is seriously affected.
 - Least-trimmed-squares essentially proceeds with OLS after eliminating the most extreme positive or negative residuals.
- LTS orders the squared residuals from smallest to largest: $(e^2)_{(1)}$, $(e^2)_{(2)}$, ..., $(e^2)_{(T)}$
- Then, LTS calculates \mathbf{b} that *minimizes the sum of only the smaller half of the residuals*.

41

Bounded Influence Regression: LTS

- LTS calculates \mathbf{b} that *minimizes the sum of only the smaller half of the residuals*:

$$\sum_{i \text{ to } m} (e^2)_{(i)} = \mathbf{0}$$

where $m = [T/2] + 1$; the square bracket indicates rounding down.

- By using only the 50% of the data that fits closest to the original OLS line, LTS completely ignores extreme outliers. The breakdown value for the LTS estimate is $(T-m)/T$.
- On the other hand, this method can misrepresent the trend in the data if it is characterized by clusters of extreme cases or if the data set is relatively small.

42

Bounded Influence Regression: LMS

- An alternative bounded influence method is *Least Median Squares (LMS)*.
- Rather than minimize the sum of the least squares function, this model minimizes the median of the squared residuals, e_i^2 .
- The breakdown value for the LTS estimate is also $(T-m)/T$.
- LMS is very robust with respect to outliers both in terms of X and Y.
- But, it performs poorly from the point of view of asymptotic efficiency. Also, relative to LMS, LTS's objective function is smoother, making the LTS estimate less jumpy -i.e., less sensitive to local effects.

43

Robust Regression: Application 1

De Long and Summers (1991) studied the national growth of 61 countries from 1960 to 1985 using OLS:

$$\text{GDP}_i = \beta_0 + \beta_1 \text{LFG}_i + \beta_2 \text{GAP}_i + \beta_3 \text{EQP}_i + \beta_4 \text{NEQ}_i + \varepsilon_i$$

where GDP growth per worker (GDP) and the regressors are labor force growth (LFG), relative GDP gap (GAP), equipment investment (EQP), and nonequipment investment (NEQ).

The REG Procedure					
Model: MCGML1					
Dependent Variable: GDP					
Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	-0.01430	0.01028	-1.39	0.1697
LFG	1	-0.02981	0.19838	-0.15	0.8811
GAP	1	0.02026	0.00917	2.21	0.0313
EQP	1	0.26538	0.06529	4.06	0.0002
NEQ	1	0.06236	0.03482	1.79	0.0787

- The OLS analysis: GAP and EQP have a significant effect on GDP at the 5% level.

44

Robust Regression: Application 1

Zaman, Rousseeuw, and Orhan (2001) used robust techniques to estimate the same model (Zambia (observation #60) an outlier):

$$GDP_i = \beta_0 + \beta_1 LFG_i + \beta_2 GAP_i + \beta_3 EQP_i + \beta_4 NEQ_i + \varepsilon_i$$

The ROBUSTREG Procedure					
Model Information					
Data Set	MFLIB.GROWTH				
Dependent Variable	GDP				
Number of Covariates	4				
Number of Observations	61				
Name of Method	M-Estimation				
Summary Statistics					
Variable	Q1	Median	Q3	Mean	Standard Deviation
LFG	0.6118	0.0239	0.02805	0.02113	0.009794
GAP	0.57955	0.8015	0.88625	0.725777	0.21807
EQP	0.0285	0.0433	0.072	0.052325	0.028622
NEQ	0.09555	0.1356	0.1812	0.13856	0.05684
GDP	0.01205	0.0231	0.03095	0.022384	0.015516
Summary Statistics					
Variable	MAD				
LFG	0.009489				
GAP	0.177764				
EQP	0.032469				
NEQ	0.052468				
GDP	0.014974				

The ROBUSTREG Procedure					
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	95% Confidence Limits	Chi-Square
Intercept	1	-0.0247	0.0097	-0.0437 -0.0058	6.53
LFG	1	0.1040	0.1867	-0.2619 0.4699	0.31
GAP	1	0.0250	0.0086	0.0080 0.0419	8.36
EQP	1	0.2968	0.0614	0.1764 0.4172	23.33
NEQ	1	0.0885	0.0328	0.0242 0.1527	7.29
Scale	1	0.0099			
Parameter Estimates					
Parameter	Pr > ChiSq				
Intercept	0.0106				
LFG	0.5775				
GAP	0.0038				
EQP	<.0001				
NEQ	0.0069				
Scale					

- Huber M-estimates: Besides GAP and EQP, the robust analysis also show NEQ has significant effect on GDP.

45

Robust Regression: Diagnostics

- It is common to analyze the residuals for outliers (as usual) and leverage points. To check for leverage points, Rousseeuw (1984) proposes a robust version of the Mahalanobis distance by using a generalized minimum covariance determinant (MCD) method.

- Mahalanobis Distance is the square root of a standard Wald distance:

$$MD(x_i) = [(x_i - \bar{x})^T \bar{C}(X)^{-1} (x_i - \bar{x})]^{1/2}$$

where \bar{x} is the mean and $\bar{C}(X)$ is the variance (scale or scatter) of X.

- Rousseeuw's Robust Distance is given by

$$RD(x_i) = [(x_i - T(X))^T C(X)^{-1} (x_i - T(X))]^{1/2}$$

where T(X) and C(X) are the robust multivariate location and scale, respectively, obtained by MCD.

46

Robust Regression: LTS - Application 1

Analysis of robust residuals. Lots of leverage observations, but only one outlier (Zambia, #60).

The ROBUSTREG Procedure

Diagnostics

Obs	Mahalanobis Distance	Robust MCD Distance	Leverage	Robust Residual	Outlier
1	2.6083	4.0639	*	-0.9424	
5	3.4351	6.7391	*	1.4200	
8	3.1876	4.6843	*	-0.1972	
9	3.6752	5.0599	*	-1.0784	
17	2.6024	3.8186	*	-1.7971	
23	2.1225	3.8238	*	1.7161	
27	2.6461	5.0336	*	0.0909	
31	2.9179	4.7140	*	0.0216	
53	2.2600	4.3193	*	-1.8082	
57	3.8701	5.4874	*	0.1448	
58	2.5953	3.9671	*	-0.0978	
59	2.9239	4.1663	*	0.3573	
60	1.8562	2.7135	*	-4.9798	*
61	1.9634	3.9128	*	-2.5959	

Diagnostics Profile

Name	Percentage	Cutoff
Outlier	0.0164	3.0000
Leverage	0.2131	3.3382

47

Robust Regression: LTS - Application 1

The analysis of robust residuals revealed Zambia (#60) as an outlier. Potentially, this can create problems for M-estimators. LTS estimation is has a better breakdown point.

The ROBUSTREG Procedure

LTS Profile

Total Number of Observations	61
Number of Squares Minimized	33
Number of Coefficients	5
Highest Possible Breakdown Value	0.4590

LTS Parameter Estimates

Parameter	DF	Estimate
Intercept	1	-0.0249
LFD	1	0.1133
GAP	1	0.0214
BQP	1	0.2669
HBQ	1	0.1110
Scale		0.0076
Wscale		0.0109

The ROBUSTREG Procedure

Diagnostics

Obs	Mahalanobis Distance	Robust MCD Distance	Leverage	Robust Residual	Outlier
1	2.6083	4.0639	*	-1.0715	
5	3.4351	6.7391	*	1.6574	
8	3.1876	4.6843	*	-0.2324	
9	3.6752	5.0599	*	-2.0896	
17	2.6024	3.8186	*	-1.6367	
23	2.1225	3.8238	*	1.7570	
27	2.6461	5.0336	*	0.2334	
31	2.9179	4.7140	*	0.0971	
53	2.2600	4.3193	*	-1.2978	
57	3.8701	5.4874	*	0.0605	
58	2.5953	3.9671	*	-0.0857	
59	2.9239	4.1663	*	0.4113	
60	1.8562	2.7135	*	-4.4984	*
61	1.9634	3.9128	*	-2.1201	

Diagnostics Profile

Name	Percentage	Cutoff
Outlier	0.0164	3.0000
Leverage	0.2131	3.3382

Resquare for LTS-estimation

Esquare 0.7417678684

Robust Regression: LTS - Application 1

After removing the outlier (Zambia), we re-estimate model:

The ROBUSTREG Procedure					
Parameter Estimates for Final Weighted LS					
Parameter	DF	Estimate	Standard Error	95% Confidence Limits	Chi-Square
Intercept	1	-0.0222	0.0093	-0.0403 -0.0041	5.75
LFG	1	0.0446	0.1755	-0.2995 0.3886	0.06
GAP	1	0.0245	0.0081	0.0085 0.0404	9.05
EQP	1	0.2824	0.0576	0.1695 0.3953	24.03
NRQ	1	0.0849	0.0311	0.0239 0.1460	7.43
Scale		0.0115			

Parameter Estimates for Final Weighted LS	
Parameter	Pr > ChiSq
Intercept	0.0165
LFG	0.7995
GAP	0.0026
EQP	<.0001
NRQ	0.0064
Scale	

Figure 7. Final Weighted LS estimates

49

Robust Regression: 3 Factor Model - Application 2

We run the 3 Fama-French factor model, for the lowest size (decile) portfolio.

```
proc robustreg data=ab;
model y1 = xm SMB HML;
output out=robout r=resid sr=stdres;
run;
```

Parameter Estimates						
Parameter	DF	Estimate	SE	95% C.I. Limits	Chi-Square	Pr > ChiSq
Intercept	1	-0.4184	0.0637	-0.5432 -0.2937	43.20	<.0001
xm	1	0.9660	0.0125	0.9415 0.9906	5947.57	<.0001
SMB	1	1.2760	0.0203	1.2363 1.3157	3969.66	<.0001
HML	1	0.4478	0.0182	0.4121 0.4835	604.18	<.0001
Scale	1	1.8128				

50

Robust Regression: Remarks

- Separated points can have a strong *influence* on statistical models
 - Unusual cases can substantially influence the fit of the OLS model. Cases that are both *outliers* and *high leverage* exert *influence* on both the slopes and intercept of the model
 - Outliers may also indicate that our model fails to capture important characteristics of the data
- Efforts should be made to remedy the problem of unusual cases before proceeding to robust regression
- If robust regression is used, careful attention must be paid to the model—different procedures can give completely different answers.

51

Robust Regression: Remarks

- No one robust regression technique is best for all data
- There are some considerations, but even these do not hold up all the time:
 - LAD regression should generally be avoided because it is less efficient than other techniques and often not very resistant
 - Bounded influence regression models, which can have a breaking point as high as 50%, often work very well with large datasets. But, they tend to perform poorly with small datasets.
- M-Estimation is typically better for small datasets, but its standard errors are not reliable for small samples. This can be overcome by using bootstrapping to obtain new estimates of the standard errors.

52

Quantile Regression

- Mosteller and Tukey (1977):

“What the regression curve does is a grand summary for the the averages of the distributions corresponding to the set of x’s. We could go further and compute several different regression curves corresponding to the various percentage points of the distribution and thus get a more complete picture.”

- One might be interested in behavior of say, lower tail of the conditional distribution rather than in its mean.
- For example, how does a 1% increase in market returns affect the returns of small size firms?

53

Quantiles: Characterizing a Distribution

- We are used to assume a distribution and describe it through its moments: mean, variance, skewness, etc. Some distributions are characterized by few parameters. For example, the normal is completely described by the mean and the variance.
- A different approach. Use quantiles instead. For example:
 - Median
 - Interquartile Range
 - Interdecile Range
 - Symmetry = $(\zeta_{.75} - \zeta_{.5}) / (\zeta_{.5} - \zeta_{.25})$
 - Tail Weight = $(\zeta_{.90} - \zeta_{.10}) / (\zeta_{.75} - \zeta_{.25})$

Quantiles

Quantiles

- We say that a firm is in the θ^b quantile if it is bigger than the proportion θ , of the reference group of firms, and smaller than the proportion $(1-\theta)$.
- The θ^b sample quantile is simply $y_{(k)}$, where k is the smallest integer such that $K/T < \theta$. (Note the relation between rank and quantile.)

55

Quantiles: Definition

Definition:

(1) Discrete RV. Given $\theta \in [0, 1]$. A θ th quantile of a discrete RV Z is any number ζ_θ such that $\Pr(Z < \zeta_\theta) \leq \theta \leq \Pr(Z \geq \zeta_\theta)$.

Example: Suppose $Z = \{3, 4, 7, 9, 9, 11, 17, 21\}$ and $\theta = 0.5$ then $\Pr(Z < 9) = 3/8 \leq 1/2 \leq \Pr(Z \geq 9) = 5/8$.

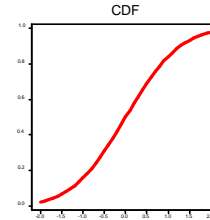
(2) Continuous RV. Let Z be a continuous r.v. with cdf $F(\cdot)$, then $\Pr(Z < z) = \Pr(Z \leq z) = F(z)$ for every z in the support and a θ th quantile is any number ζ_θ such that $F(\zeta_\theta) = \theta$

- If F is continuous and strictly increasing then the inverse exists and $\zeta_\theta = F^{-1}(\theta)$.

Quantiles: CDF and Quantile Function

- Cumulative Distribution Function

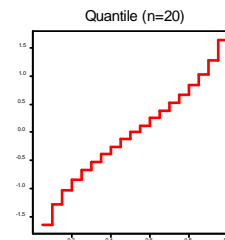
$$F(y) = \text{Prob}(Y \leq y)$$



- Quantile Function

$$Q(\theta) = \min(y : F(y) \leq \theta)$$

=> Discrete step function



Quantiles

- It can be shown that quantile (θ) is the solution to

$$\arg \min_{\xi} \frac{1}{T} \left\{ \sum_{y_t \geq \xi} \theta |y_t - \xi| + \sum_{y_t < \xi} (1 - \theta) |y_t - \xi| \right\}$$

- If $\theta = 1/2$, then this becomes $\arg \min_{\mu} \frac{1}{T} \sum_{t=1}^T |y_t - \xi|$, which yields a f.o.c.:

$$0 = (-1/T) \sum_t [\text{sgn}(y_t - \xi)]$$

where sng (“*signum*”) function: $\text{sgn}(u) = 1 - 2 I[u < 0]$, (defined to be right-continuous).

=> the sample median, $\zeta_{\theta=0.50}$, solves this problem (easier to visualize with expectations).

Quantile Regression

- Basset and Koenker (1978, JASA) suggest simply replacing the ζ in the definition of the quantile estimator

$$\arg \min_{\xi} \sum_{y_t \geq \xi} \theta |y_t - \xi| + \sum_{y_t < \xi} (1 - \theta) |y_t - \xi|$$

with $\mathbf{X}_t' \boldsymbol{\beta}$ to get the *quantile regression*

$$\arg \min_{\boldsymbol{\beta}} \sum_{y_t \geq \mathbf{X}_t' \boldsymbol{\beta}} \theta |y_t - \mathbf{X}_t' \boldsymbol{\beta}| + \sum_{y_t < \mathbf{X}_t' \boldsymbol{\beta}} (1 - \theta) |y_t - \mathbf{X}_t' \boldsymbol{\beta}| = \sum_{\varepsilon_t \geq 0} \theta |\varepsilon_t| + \sum_{\varepsilon_t < 0} (1 - \theta) |\varepsilon_t|$$

- If $\theta = 1/2$, then this becomes LAD estimation. We have a symmetric weighting of observations with positive and negative residuals. But, if $\theta \neq 1/2$, the weighting is asymmetric.

59

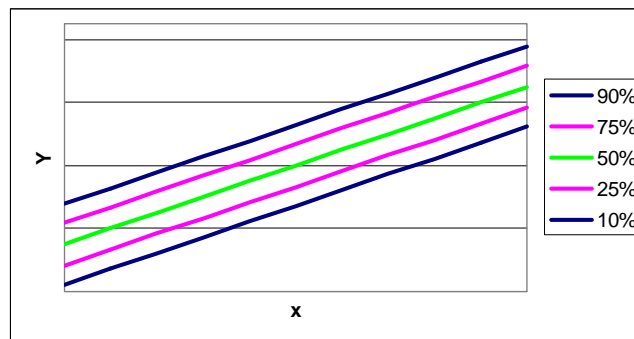
Quantile Regression

- We define a family of regressions:

$$\zeta_{\theta} = Q(y | \mathbf{x}, \theta) = \mathbf{X}' \boldsymbol{\beta}_{\theta}, \quad \theta \in [0, 1]$$

- Median regression is obtained by setting $\theta = .50$:

$$\zeta_{\theta=.50} = Q(y | \mathbf{x}, .50) = \mathbf{X}' \boldsymbol{\beta}_{\theta=.50}$$



60

Quantile Regression

Note: Median regression estimated by LAD. It estimates the same parameters as OLS if symmetric conditional distribution.

- We assume *correct specification* of the quantile, $Q(y | \mathbf{x}, \theta) = \mathbf{X}'\boldsymbol{\beta}_\theta$. That is, $\mathbf{X}'\boldsymbol{\beta}$ is a particular linear combination of the independent variables such that

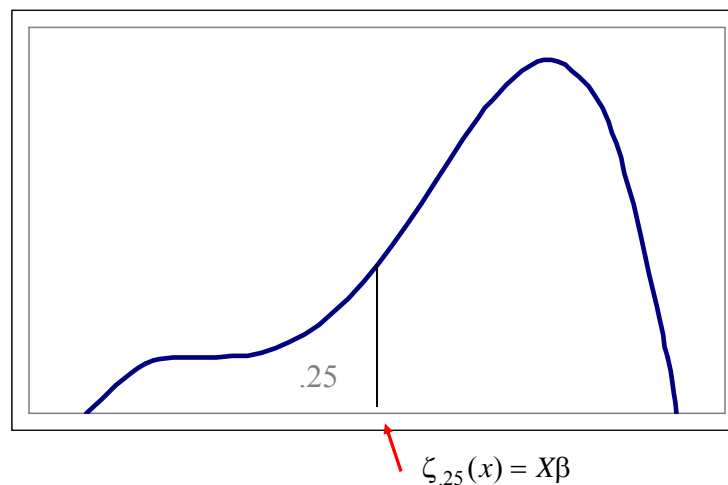
$$\theta = \Pr(Y \leq \zeta_\theta(X) | X) = \Pr(Y \leq X\boldsymbol{\beta}) = F(\zeta_\theta(X) | X)$$

Q: Why use quantile (median) regression?

- Semiparametric
- Robust to some extensions (heteroscedasticity?)
- Complete characterization of conditional distribution.

61

Quantile Regression



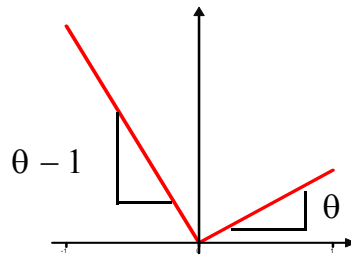
Quantile Regression: Loss Function

- Different from LS, now we minimize an asymmetric absolute loss function, given by

$$\arg \min_{\beta} \rho_{\theta}(y_t, X_t' \beta) = \arg \min_{\beta} \sum_{y_t \geq X_t' \beta} \theta |y_t - X_t' \beta| + \sum_{y_t < X_t' \beta} (1 - \theta) |y_t - X_t' \beta|$$

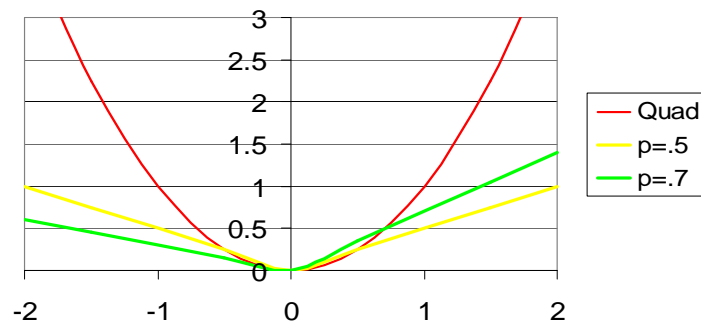
for some θ .

- We call ρ_{θ} the *tilted absolute value function*. It is convex. The local minimum is a global one, which assures uniqueness (and identification).



Quantile Regression: Loss Function

Absolute Loss vs. Quadratic Loss over errors



A quadratic loss penalizes large errors very heavily. When $p=.5$ our best predictor is the median; it does not give as much weight to outliers. When $p=.7$ the loss is asymmetric; large positive errors are more heavily penalized than negative errors.

Quantile Regression: Estimation

- Optimization problem:

$$\min_{\beta} \sum_{\varepsilon_t \geq 0} \theta |\varepsilon_t| + \sum_{\varepsilon_t < 0} (1 - \theta) |\varepsilon_t| = \sum_{t=1}^T (\theta - I[y_t < \theta]) \varepsilon_t$$

- Simple intuition: number of negative residuals $\leq T\theta \leq$ number of negative residuals + number of zero residuals.

- Since the loss function is piecewise linear, solving it is a linear programming problem. Trick: replace absolute values by positivity constraints. That is,

$$\min_{\beta} \left\{ \sum_{t=1}^T \theta \varepsilon_t^+ + (1 - \theta) \varepsilon_t^- = \theta \mathbf{1}' \varepsilon^+ + (1 - \theta) \mathbf{1}' \varepsilon^- \right\}$$

$$s.t. \quad y = X\beta + \varepsilon^+ - \varepsilon^- \quad (\varepsilon_t^- \leq y_t - X_t\beta \leq \varepsilon_t^+)$$

$$\varepsilon_t^+ \geq 0, \quad \varepsilon_t^- \geq 0$$

65

Quantile Regression: Estimation

- The usual software packages will use the Barrodale and Roberts (1974) simplex algorithm or a Frisch-Newton (FN) algorithm.
- For large data sets, the FN method is used. It combines a log-barrier Lagrangian (Frisch part) with steepest descent steps (Newton part). For very large data sets, FN algorithm is combined with a preprocessing step, which makes the computations faster.
- Solution at vertex of feasible region. The solution need not be unique (along the edge). The fitted line will go through k data points.
- Well known program in R, written by Koenker and described in Koenker's Vignette article (2005).

66

Quantile Regression: Optimality

- Proposition

Under the asymmetric absolute loss function Q_θ a best predictor of Y given $X=x$ is the θ^{th} conditional quantile, ζ_θ .

Example: Let $\theta = .5$. Then, the best predictor is the median fitted value.

- That is, under asymmetric absolute loss, the quantile regression estimator is more efficient than OLS.

- We offer this without proof. The proof would be similar in construction to the Gauss-Markov Theorem, which states that the conditional mean is best linear unbiased.

Properties of the Estimator

- Consistency

Consistency of $\hat{\beta}_\theta$ is easy. The minimand $S_n(\cdot)$ is continuous in β with probability 1. In fact, $S_n(\cdot)$ is convex in β ; then, consistency follows if S_n can be shown to converge *pointwise* to a function that is uniquely minimized at the true value β_θ .

- To prove consistency, we impose conditions on the model:

1. The data $(x_i; y_i)$ are *i.i.d.* across i
2. The regressors have bounded second moment.
3. $\epsilon_i | X_i$ is continuously distributed; with conditional density $f_\epsilon(\epsilon_i | X_i)$ satisfying the conditional quantile restriction.
4. The regressors and error density satisfy a local identification condition: $C = E[f'_\epsilon(0) \mathbf{x}\mathbf{x}']$ is a pd matrix.

Properties of the Estimator

- Asymptotic Normality (under *i.i.d* assumption)

The lack of continuously differentiable $S_n(\beta)$ complicates the usual derivation of asymptotic normality (through Taylor's expansion).

- But, an approximate f.o.c. can be used -through $\text{sgn}(\cdot)$. Additional conditions (*stochastic equicontinuity*) need to be established before using the Lindeberg-Levy CLT, which establishes:

$$\sqrt{T}(\hat{\beta}_\theta - \beta_\theta) \xrightarrow{L} N(0, \Lambda_\theta)$$

where

$$\Lambda_\theta = \theta(1-\theta) (E[f_\varepsilon(0 | x_i)x_i x_i'])^{-1} E[x_i x_i'] (E[f_\varepsilon(0 | x_i)x_i x_i'])^{-1}$$

- We have a sandwich estimator. The variance matrix depends on the unknown $f_\varepsilon(\cdot | \mathbf{x})$ and the X, at which the covariance is being evaluated.

Properties of the Estimator

- We need to estimate $E[f_\varepsilon(0 | \mathbf{x}) \mathbf{x}\mathbf{x}']$, complicated without knowing $f_\varepsilon(\cdot | \mathbf{x})$! It can be done through non-parametric kernel estimation.

- When the error is independent of \mathbf{x} -i.e., $f_\varepsilon(\varepsilon_i | X_i) = f_\varepsilon(\varepsilon_i)$ -, then the coefficient covariance reduces to

$$\Lambda_\theta = \frac{\theta(1-\theta)}{f_\varepsilon^2(0)} (E(\mathbf{x}\mathbf{x}'))^{-1}$$

where

$$\hat{E}(\mathbf{x}\mathbf{x}') = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$$

- The variance is related to a Bernoulli variance $[\theta(1-\theta)]$ -divided by the square density of Y at the quantile, analagous to a sample size.

Properties of the Estimator

- The previous results can be extended to multivariate cases –i.e., joint estimates of several quantiles. We obtain convergence to a multivariate normal distribution.
- In general, the quantile regression estimator is more efficient than OLS. But, efficiency requires knowledge of the true error's pdf.
- Robust to outliers. As long as the sign of the residual does not change, any y_i can be arbitrarily changed without shifting the conditional quantile line.
- The regression quantiles are correlated.

Partial Effects and Prediction

- The marginal change in the Θ th conditional quantile due to a marginal change in the j th element of \mathbf{x} .

$$\frac{\partial Q_{\theta}(y_i | X_i)}{\partial x_{i,j}}$$

- Under linearity, the effect will be β_j . But, if non-linearities are included, the partial effect will be a function of \mathbf{x} .

Note: There is no guarantee that the i th observation will remain in the same quantile after $x_{i,j}$ changes.

- Using $\hat{\beta}_{\theta}$ and \mathbf{X} values, predicted values of \hat{y}_{θ} , can be computed. Suppose we have $\mathbf{X}=\mathbf{x}_0$, the predicted 90th quantile is $\mathbf{x}_0' \hat{\beta}_{.90}$.

Hypothesis Testing: Standard Errors

- Given asymptotic normality, one can construct asymptotic t-statistics for the coefficients. But which standard errors should be used?
- We can use the asymptotic estimator, but in non-*i.i.d.* situations is complicated. Inversion of a rank test --Koenker (1994, 1996)-- can be used to construct C.I.'s in a non-*i.i.d.* error context.
- Bootstrapping works well. Parzen, Wei, and Ying (1994) have suggested that rather than bootstrapping $(x_i; y_i)$ pairs, instead bootstrap the quantile regression gradient condition. It produces a pivotal approach.

Hypothesis Testing

- Alternatively, confidence regions for the quantile regression parameters can be computed from the empirical distribution of the sample of bootstrapped $\mathbf{b}_j(\theta)$'s, the so-called percentile method.
- These procedures can be extended to deal with the joint distribution of several quantile regression estimators $\{\mathbf{b}_j(\theta_k), k = 1; 2, \dots, K\}$. This would be needed to test equality of slope parameters across quantiles.
- The error term may be heteroscedastic. Efficiency issue. There are many tests for heteroscedasticity in this context.
- A test for symmetry, resembling a Wald Test, can be constructed which could not be done under Least Squares estimation.

Crossings

- Since quantile regressions are typically estimated individually, the quantile curves can cross, leading to strange (an invalid) results.
- Crossings problems increase with the number of regressors..
- Simultaneous estimation, with constraints are one solution.
- Individual specification of each quantile also works. For example:

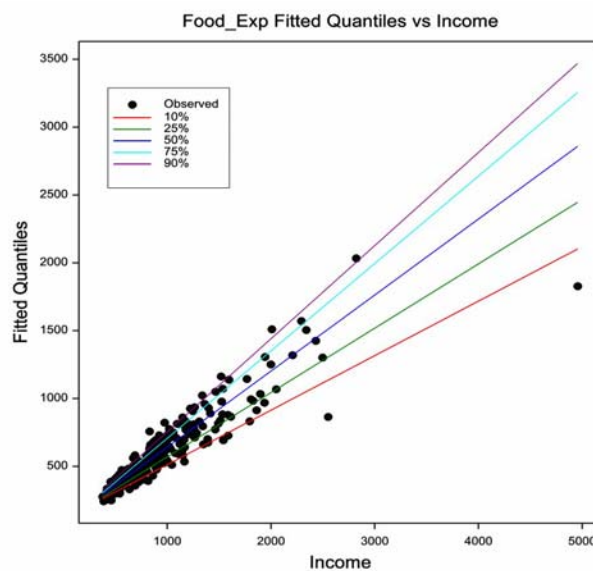
$$y = X\beta_0 + \varepsilon^0, \quad P[\varepsilon^0 < 0 | X] = \theta_0 \quad (\text{say, } \theta_0 = .5)$$

$$y = X\beta_0 - \exp(X\beta_1) + \varepsilon^1, \quad P[\varepsilon^1 < 0 | X] = \theta_1 \quad (\text{say, } \theta_0 = .25)$$

$$y = X\beta_0 + \exp(X\beta_2) + \varepsilon^2, \quad P[\varepsilon^2 < 0 | X] = \theta_2 \quad (\text{say, } \theta_0 = .75)$$

Note: Since $\exp(\cdot)$ is positive, the quantiles by design never cross.

Quantile Linear Regression: Application 1

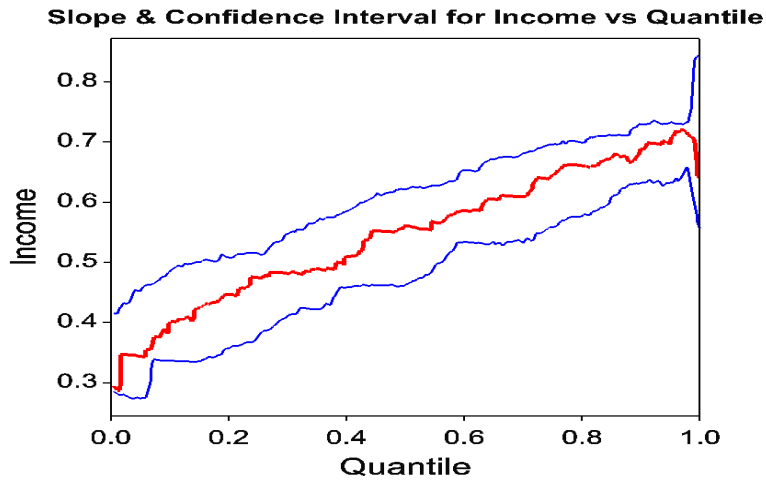


Food Expenditure vs Income

Engel 1857 survey of 235 Belgian households

Q: Change of slope at different quantiles?

Quantile Linear Regression: Application 1



Note: Variation of Parameter with Quantiles.

Quantile Linear Regression: Application 2

```
Quantile Regression Model. Quantile = .250000
Linear Programming estimation method
LHS=HHNINC  Mean = 44583
            Standard deviation = 21550
            Number of observs. = 3377
            Minimum = 04000
            t= .25000 quantile = 30000
            Maximum = 3.00000
Model size  Parameters = 5
            Degrees of freedom = 3372
Residuals  Sum of squares = 193.75951
            Standard error of e = 20226
Fit         R-squared = .12721
            PseudoR2=1-F(0)/F(b) = .11046
Not using OLS or no constant. Rsquared may be <= 0
Functions F= Sum r(t)[y(i)-x(i)b] = 164.31749
            F0=Sum r(t)[y(i)-Qy(t)] = 184.72281
            r(t)[u]=t*u-u*[u<0]. t= .250000
```

HHNINC	Coefficient	Standard Error	z	Prob. z >Z*	95% Confidence Interval	
Constant	-.07580***	.01839	-4.12	.0000	-.11185	-.03975
AGE	-.00036**	.00016	-2.25	.0244	-.00068	-.00005
EDUC	.02393***	.00137	17.51	.0000	.02125	.02661
MARRIED	.11459***	.00547	20.96	.0000	.10398	.12531
HSAT	.00773***	.00122	6.31	.0000	.00533	.01013
Constant	-.01504	.03479	-.43	.6656	-.08323	.05315
AGE	-.00035	.00039	-.90	.3669	-.00112	.00041
EDUC	.02707***	.00167	16.19	.0000	.02379	.03035
MARRIED	.11361***	.01115	10.19	.0000	.09175	.13547
HSAT	.00777***	.00195	3.99	.0001	.00396	.01158
Constant	.03738	.03246	1.15	.2495	-.02624	.10099
AGE	.00020	.00039	.51	.6100	-.00057	.00097
EDUC	.03240***	.00237	13.68	.0000	.02776	.03704
MARRIED	.08042***	.01112	7.23	.0000	.05863	.10222
HSAT	.00693***	.00231	3.00	.0027	.00240	.01145

Note: ***, **, * ==> Significance at 1%, 5%, 10% level.

$\alpha = .25$

$\alpha = .50$

$\alpha = .75$

Quantile Linear Regression: SAS - Application 3

```
proc quantreg data=ab ;
model y1 = xm SMB HML /quantile=0.25 0.5 0.75
run;
```

The QUANTREG Procedure
Quantile and Objective Function

Quantile	0.25				
Parameter	DF	Estimate	95% Confidence		Limits
Intercept	1	-1.6310	-1.7793	-1.5164	
xm	1	0.9855	0.9477	1.0069	
SMB	1	1.2018	1.1505	1.3219	
HML	1	0.5071	0.4250	0.5615	

Quantile	0.75				
Parameter	DF	Estimate	95% Confidence		Limits
Intercept	1	0.9056	0.6957	1.1344	
xm	1	0.9919	0.9626	1.0535	
SMB	1	1.4267	1.3340	1.5025	
HML	1	0.5435	0.4593	0.6213	

Heteroscedasticity

- Model: $y_i = x_i'\beta + \varepsilon_i$, with *i.i.d.* errors.
 - The quantiles are a vertical shift of one another.
- Model: $y_i = x_i'\beta + \sigma(x_i) \varepsilon_i$, errors are now heteroscedastic.
 - The quantiles now exhibit a location shift as well as a scale shift.
- Khmaladze-Koenker Test Statistic

Quantile Regression: Bibliography

- Buchinsky, M. (1994), “Changes in the u.s. wage structure 1963-1987: Application of quantile regression,” *Econometrica*, 62, 405-458.
- Koenker and Hulloch (2001), “Quantile Regression,” *Journal of Economic Perspectives*, Vol. 15, Pps. 143-156.
- Koenker (2005), **Quantile Regression**, Cambridge University Press.

Quantile Regression

- S+ Programs - [Lib.stat.cmu.edu/s](http://lib.stat.cmu.edu/s)
- www.econ.uiuc.edu/~roger
- [http://Lib.stat.cmu.edu/R/CRAN](http://lib.stat.cmu.edu/R/CRAN)
- TSP
- Limdep