# Outlier detection and treatment

LECTURE 12

C4D2 TRAINING

1

---

Today is mainly about outliers

1) Definitions
   What do we mean by an outlier, exactly?

2) Motivation
   Do outliers really matter?

3) Detection
   How to detect outliers?

4) Treatment
   How to deal with outliers?

C4D2 TRAINING

2

---

# Definitions

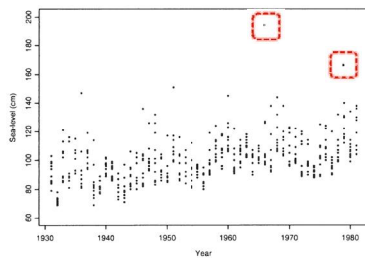C4D2 TRAINING

3

## What is an outlier?

- An outlier is an observation "that appears to deviate markedly from other members of the sample in which it occurs" (Grubbs, 1969)

- Note: we focus on univariate outliers, those found when looking at a distribution of values in a single dimension (e.g. income).

*Outliers in Statistical Data*

VIC. BARNETT
*University of Sheffield*

and

TOBY LEWIS
*University of Hull*

*John Wiley & Sons*
*Chichester · New York · Brisbane · Toronto*

C4D2 TRAINING

---

## Highest sea-levels in Venice



C4D2 TRAINING
5

---

## Other classical definitions

- An outlier is "an observation that deviates so much from other observations as to arouse suspicion that it was generated by a different mechanism" (Hawkins 1980)

- Aguinis et al (2013) provide 14 definitions of outliers based on a litterature review of 28 papers.

C4D2 TRAINING
6

## What causes outliers?

- Human errors, e.g. data entry errors

- Instrument errors, e.g. measurement errors

- Data processing errors, e.g. data manipulation

- Sampling errors, e.g. extracting data from wrong sources

- Not an error, the value is extreme, just a 'novelty' in the data

## A dilemma

- Outliers can be genuine values

- The trade-off is between the loss of accuracy if we throw away "good" observations, and the bias of our estimates if we keep "bad" ones

- The challenge is twofold:

    1. to figure out whether an extreme value is good (genuine) or bad (error)

    2. to assess its impact on the statistics of interest
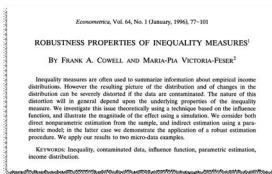
## Do outliers matter?

## Theory first

- Three papers:

  I. 1996a
     Frank Cowell and Maria-Pia Victoria-Feser

  II. 2007
      Frank Cowell and Emmanuel Flachaire (*)

  III. 1996b
       Frank Cowell and Maria-Pia Victoria-Feser

---

## Outliers and inequality measures – I
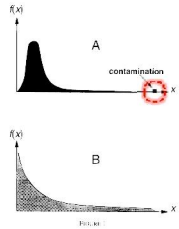Cowell and Victoria-Feser (1996a)

- This is a beautiful paper

- Explains why outliers (contaminants) are a serious threat to most inequality measures.

- "if the mean has to be estimated from the sample then all scale independent or translation independent and decomposable measures have an unbounded influence function" (p. 89)

- An unbounded IF is a catastrophe.

---

## *The influence function

- $F$ — Ideal data, no contaminants
- $Gini_{TRUE} = I(F)$ — "true" Gini index

- $G = (1 - \delta)F + \delta H$ — Real-world data, with $\delta$% contaminants
  $0 \leq \delta \leq 1$

- $Gini_{ESTIMATED} = I(G)$ — estimated Gini index

- The influence function, IF: $IF = \lim_{\delta \to 0} \frac{I(G) - I(F)}{\delta}$

## The catatrophe



- Suppose the shape of the income distribution is represented by the continuous frequency distribution in part A

- Suppose that in the sample there are some rogue observations represented by the point mass labelled "contamination".

- Then, according to inequality statistics that are sensitive to the top end of the distribution, the income distribution in A will be indistinguishable from that represented in B (that is, IF is unbounded).

---

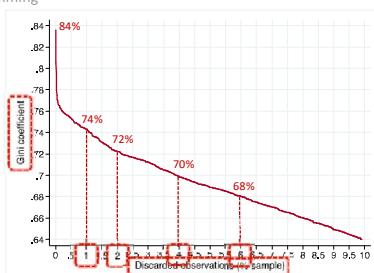## In practice
Hlasny and Verme (2018: 191)

- Many researchers routinely trim outliers or problematic observations or apply top coding with little consideration of the implications for the measurement of inequality

- One example to illustrate

14

---

## Sensitivity of the Gini index to extreme values
iterative trimming

## *Outliers and inequality measures – II
Cowell and Flachaire (2007)

Abstract

We examine the statistical performance of inequality indices in the presence of extreme values in the data and show that these indices are very sensitive to the properties of the income distribution. Estimation and inference can be dramatically affected, especially when the tail of the income distribution is heavy, even when standard bootstrap methods are employed. However, use of appropriate semiparametric methods for modelling the upper tail can greatly improve the performance of even those inequality indices that are normally considered particularly sensitive to extreme values.
© 2007 Published by Elsevier B.V.

JEL classification: C1; D63

Keywords: Inequality measures; Statistical performance; Robustness

- Explains how and why outliers are a serious threat to most inequality measures.

C4D2 TRAINING

---

## *How rapidly the catastrophe occurs
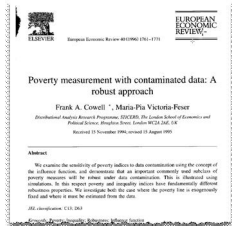Rates of increase to infinity of the influence function

- Let us concentrate only on the extremes of the income distribution. Data contamination can occur at very high incomes (say at a point $z$ that approaches infinity) or at very low incomes ($z \approx 0$).

| Measure | Generalised entropy, $I_E$ | | | | Atkinson, $I_A$ | | | LogVar | Gini |
|---|---|---|---|---|---|---|---|---|---|
| | $\alpha > 1$ | $0 < \alpha \leq 1$ | $\alpha = 0$ | $\alpha < 0$ | $0 < \varepsilon < 1$ | $\varepsilon = 1$ | $\varepsilon > 1$ | | |
| $z \to \infty$ | $z^\alpha$ | $z$ | $z$ | – | $z$ | $z$ | – | $z$ | $z$ |
| $z \to 0$ | – | – | $\log z$ | $z^\alpha$ | – | $\log z$ | $z^{1-\varepsilon}$ | $(\log z)^2$ | – |

- Result 1: GE measures with $\alpha > 1$ are very sensitive to high incomes in the data.

- Result 2: GE measures with $\alpha < 0$, and Atkinson measures with $\varepsilon > 1$ are very sensitive to small incomes in the data.

- We will return on this catastrophe in due time, later during this workshop.

C4D2 TRAINING

---

## Outliers and poverty measures
Cowell and Victoria-Feser (1996b)

Abstract

We examine the sensitivity of poverty indices to data contamination using the concept of the influence function, and demonstrate that an important commonly used subclass of poverty measures will be robust under data contamination. This is illustrated using simulations. In this respect poverty and inequality indices have fundamentally different robustness properties. We investigate both the case where the poverty line is exogenously fixed and where it must be estimated from the data.

JEL classification: C13; D63

- Explains why outliers only rarely are a serious threat to most poverty measures.

- Poverty measures are not sensitive to the values (real or contaminated) of the incomes of the rich

C4D2 TRAINING

Recap

- The answer to the question on whether outliers matter depends on the statistic of interest

- Inequality: both theory (unbounded IF) and practice (incremental truncation) suggest that they matter (tremendously). Not taking this issue into proper account puts inequality comparisons at risk.
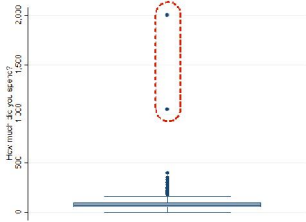
- Poverty: not so much

# How to detect outliers?

Visual inspection

- Our procedures are part graphical, and part automatic. For each commodity, we draw histograms and one-way plots of the logarithms of the unit values, using each to detect the presence of gross outliers for further investigations. [...] [Automatic method] does not remove the need for the graphical inspection
(Deaton and Tarozzi 2005)

## Visual inspection
Malawi IHS3, Cassava tuber expenditure



Malawi - Integrated Household Panel Survey 2010-2013 (Short-Term Panel, 204 EAs)

---

## Visual inspection
Malawi IHS3, Cassava tuber expenditure

- Example 1: look at descriptive statistics

---

## Visual inspection
Malawi IHS3, Cassava tuber expenditure

- Example 2: graph the distribution of the data

## Visual inspection
Malawi IHS3, Cassava tuber expenditure

- Example 3: use graphical diagnostic tools, e.g. the boxplot graph

---

## Statistical methods

- The literature is rich with methods to identify outliers; in practice, most methods used in empirical work hinge on the underlying distribution of the data.

- The idea is simple:
  - transform the variable to induce normality
  - set thresholds to identify extreme values

---

## Transform the variable to induce normality

- The easiest transformation relies on taking the logarithm of the variable of interest
- The log "squeezes" large values more, so that skewed distributions become more symmetrical and closer to a Normal distribution.

## Set a threshold

- We must specify a threshold for deciding whether each observation is 'too extreme' (outlier or not?)

- Common 'thumb-rule' thresholds : an observation is considered an outlier if it is more than 2.5, 3, 3.5 standard deviations far from the mean of the distribution

- In formulas: $x$ is an outlier if $\quad x > \bar{x} + z_\alpha\, s$

  where $z_\alpha$ equals, say, 2.5.

- We can express the same criterion as $\quad \dfrac{x - \bar{x}}{s} > z_\alpha$

  where the left-hand side is called a z-score (a variable with mean = 0 and var = 1)

---

## Why 2.5, 3, or any other number?

- Under the assumption of normality:



- $z_\alpha = 2.5$ implies that outliers are in the region where $\alpha = 0.5$ percent of other observations normally are.

---

## Deaton and Tarozzi (2005)



In the case of India, D&T (2005) flagged as outliers prices whose logarithms exceeded the mean of logarithms by more than 2.5 standard deviations:
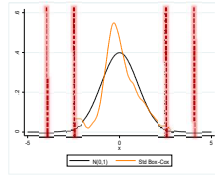
$$\frac{ln(x) - E[ln(x)]}{sd[ln(x)]} > 2.5$$

## Transformation and thresholds

Raw untransformed data

Transformed data

---

## Two questions

1) How good is such an approach?

2) What to do after flagging outliers?

---

## How good is such an approach?

- Log-transformation is very basic – how to deal with negative values?

- Not recommended when the log-distribution can not be assumed to be a Normal distribution

- Why should we set the threshold using the mean and standard deviation, which are sensitive to extreme values, if this is exactly what we are worried about?

$$\frac{ln(x) - E[ln(x)]}{sd[ln(x)]} > 2.5$$

- We can do better

## *The Box-Cox transformation

Building a household consumption database for
the calculation of poverty PPPs

Technical note

DRAFT 1.0

Olivier Dupriez, World Bank
March 2007

- The Box-Cox transformation:

$$y_h^{(\lambda)} = \begin{cases} (y_h^\lambda - 1)/\lambda & \text{if } \lambda \neq 0 \\ \ln y_h & \text{if } \lambda = 0 \end{cases}$$

- The transformed variable is often remarkably close to a Normal dn.

- Outliers are identified if:
  $y_h > $ 75th percentile + 5 × IQR

C4D2 TRAINING

---

## *The Box-Cox method: An assessment

- Only applies to strictly positive variables (e.g., it does not necessarily work with income)

- Calculation is cumbersome, and often problematic

C4D2 TRAINING                                                    35

---

## A popular strategy
robustification

- While there is no agreement on the best method, a common solution is to use robust measures of scale and location to set the threshold for flagging outliers

- the idea is to replace the sample average $\bar{x}$ with a robust estimator (e.g. the median), and the standard deviation $s$ with a robust estimator. A popular option is the median absolute deviation (MAD).

C4D2 TRAINING                                                    36

## The median absolute deviation (MAD)

$$z_h = \frac{x_h - \bar{x}}{s} \qquad z_h = \left| \frac{x_h - med[x_h]}{MAD} \right|$$

$$MAD = b \times med|x - med[x]|$$

b = 1.4826

if the distribution is Gaussian

---

## We can do better
Rousseeuw and Croux (1993, JASA)

### Alternatives to the Median Absolute Deviation

Peter J. Rousseeuw and Christophe Croux*

In robust estimation one frequently needs an initial or auxiliary estimate of scale. For this one usually takes the median absolute deviation $MAD_n = 1.4826$ $med_i\{|x_i - med_j x_j|\}$, because it has a simple explicit formula, needs little computation time, and is very robust as witnessed by its bounded influence function and its 50% breakdown point. But there is still room for improvement in two areas: the fact that $MAD_n$ is aimed at symmetric distributions and its low (37%) Gaussian efficiency. In this article we set out to construct explicit and 50% breakdown scale estimators that are more efficient. We consider the estimator $S_n = 1.1926$ $med_i\{med_j|x_i - x_j|\}$ and the estimator $Q_n$ given by the .25 quantile of the distances $\{|x_i - x_j|; i < j\}$. Note that $S_n$ and $Q_n$ do not need any location estimate. Both $S_n$ and $Q_n$ can be computed using $O(n \log n)$ time and $O(n)$ storage. The Gaussian efficiency of $S_n$ is 58%, whereas $Q_n$ attains 82%. We study $S_n$ and $Q_n$ by means of their influence functions, their bias curves (for implosion as well as explosion), and their finite-sample performance. Their behavior is also compared at non-Gaussian models, including the negative exponential model where $S_n$ has a lower gross-error sensitivity than the MAD.

KEY WORDS: Bias curve; Breakdown point; Influence function; Robustness; Scale estimation.

---

## Rousseeuw and Croux (1993)

- Rousseeuw and Croux (1993) propose to substitute the MAD with a different estimator:

- $S = c \times med_i\{med_j|x_j - x_i|\}$

- For each i we compute the median of $|x_i - x_j|$ (j = 1, ..., n ). This yields n numbers, the median of which gives our final estimate S.

$$z_h = \left| \frac{x_h - med[x_h]}{S} \right| \qquad c = 1.1926 \text{ at the Gaussian model.}$$

Recap

- "take the log and run" is not a recommended practice
- taking the log and robustifying the z-score is a better practice
- Belotti and Vecchi (2019) provide `outdetect.ado`

---

Malawi, 2013

| Method | Overall | Left | Right |
|--------|---------|------|-------|
| Z-score | 2.08 | 0.20 | 1.88 |
| MAD-score | 3.05 | 0.35 | 2.70 |
| S-score | 3.02 | 0.35 | 2.67 |
| Q-score | 3.00 | 0.35 | 2.65 |

- 'take the log and run': 2.08% of outliers (most of which in the right tail)
- 'take the log, robustify the z-score, and run': 3.00% (most of which in the right tail)

---

# How to deal with outliers?
(in one slide)

## Treatment of outliers

Three main methods of dealing with outliers, apart from removing them from the dataset:

1) reducing the weights of outliers (trimming weight)
2) changing the values of outliers (Winsorisation, trimming, imputation)
3) using robust estimation techniques (M-estimation).

- Documentation, transparency & reproducibility

C4D2 TRAINING

---

## Lessons learned

- Outliers can be genuine observations… be gentle to the data and document each and every step of the data processing

- As far as inequality is concerned, outliers are the worst enemy (unbounded IF)

- Outlier detection:

  - go beyond the "take the log and run" strategy. It works well only if you can describe the data with a Gaussian distribution. Typically, however, distributions are skewed.

  - Use a "take the log, robustify the z-score and run", strategy.

- Outlier treatment: it depends. Quantile regression is a good candidate.

C4D2 TRAINING                                                                44

---

## References

**Required readings**

Barnett, V., & Lewis, T. (1994). Outliers in Statistical Data. 3rd edition. J. Wiley & Sons (Chapter 1 & 2)

**Suggested readings**

Alvarez, E., Garcıa-Fernández, R. M., Blanco-Encomienda, F. J., & Munoz, J. F. (2014). The effect of outliers on the economic and social survey on income and living conditions. World Acad. Sci., Eng. Technol., Int. J. Soc., Behav., Educ., Econ., Bus. Ind. Eng, 8, 3276-3280.

Belotti, F., & Vecchi, G. (2019). Take the Log and Run: Outliers and Welfare Measurement, mimeo.

Cowell, F. A., & Flachaire, E. (2007). Income distribution and inequality measurement: The problem of extreme values. Journal of Econometrics, 141(2), 1044-1072.

Cowell, F., & Victoria-Feser, M. (1996). Robustness Properties of Inequality Measures. Econometrica, 64(1), 77-101.

Cowell, F. A., & Victoria-Feser, M. P. (1996). Poverty measurement with contaminated data: A robust approach. European Economic Review, 40(9), 1761-1771.

Deaton, A., & Tarozzi, A. (2005). "Prices and Poverty in India." The Great Indian Poverty Debate. New Delhi : MacMillan.

Dupriez, O. (2007). Building a household consumption database for the calculation of poverty PPPs. Technical note. Available at: http://go.worldbank.org/4YG7I5RGT0.

Grubbs, F. E. (1969). Procedures for detecting outlying observations in samples. Technometrics, 11(1), 1-21.

Hlasny, V., & Verme, P. (2018). Top Incomes and Inequality Measurement: A Comparative Analysis of Correction Methods Using the EU SILC Data. Econometrics, 6(2), 30.

Mancini, G., & Vecchi, G. (2019). On the Construction of a Welfare Indicator for Inequality and Poverty Analysis, mimeo.

OECD (2013). OECD Guidelines for Micro Statistics on Household Wealth

Rousseeuw, P. J., & Croux, C. (1993). Alternatives to the median absolute deviation. Journal of the American Statistical association, 88(424), 1273-1283.

C4D2 TRAINING

Thank you for your attention

Homework

Exercise 1 - Engaging with the literature

World Academy of Science, Engineering and Technology
International Journal of Economics and Management Engineering
Vol:8, No:10, 2014

The Effect of Outliers on the Economic and Social
Survey on Income and Living Conditions

Encarnación Álvarez, Rosa M. García-Fernández, Francisco J. Blanco-Encomienda, Juan F. Muñoz

Summarize the main conclusions
of the paper: do outliers matter?
Why or why not?

Exercise 2 - Do-it-yourself….

| English | Stata/R/SPSS/Excel/… |
|---|---|

1) Generate a log-normal looking wealth distribution

2) Estimate the Gini index

3) Contaminate the distribution with a few extreme values
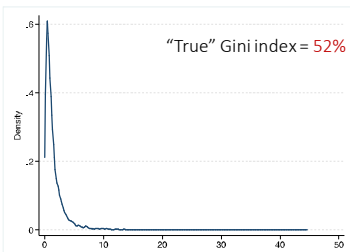
4) Re-estimate the Gini index

```
clear

set obs 5000
set seed 190607
gen n  = rnormal(0,1)
gen ln = exp(n)

* simulate order of magnitude mistake:
* take 100 obs around the median
* of the distribution and multiply
* them by 100

sort ln

gen cont100 = 1
replace cont100 = 100 in 2480/2520
gen ln_cont100 = ln*cont100
```
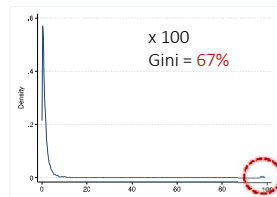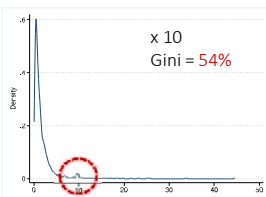
C4D2 TRAINING

---

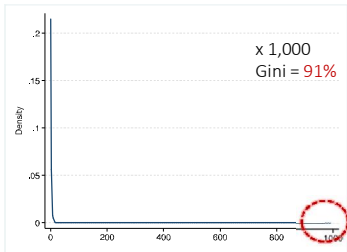"True" distribution



"True" Gini index = 52%

C4D2 TRAINING

---

Contamination

40 out of 5,000 observations (less than 1%) are "contaminated"



x 10
Gini = 54%

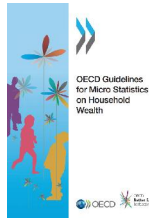x 100
Gini = 67%

C4D2 TRAINING

## Contamination

x 1,000
Gini = 91%

## Exercise 3 – Inequality measures

- Comment on table 7.3 from OECD (2013) p.172 (see next slide).
- What can you say about the sensitivity of estimates to the treatment of outliers?

OECD Guidelines for Micro Statistics on Household Wealth

53

## Exercise 3 – Inequality measures
OECD (2013)

Table 7.3. **Effect of the treatment of outliers on summary measures of wealth inequality in the United States, 2007**

|  | Raw | Shave top and bottom 1% | Shave top 1% and bottom 0.5% |
|---|---|---|---|
| Mean | 556 846 | 378 215 | 559 361 |
| Median | 120 780 | 120 780 | 123 800 |
| Gini | 0.82 | 0.74 | 0.81 |
| ½CV² | 18.1 | 2.4 | 14.6 |
| P90/P10 | 30 000 | 3 369 | 3 061 |
| P75/P25 | 26.3 | 24.5 | 24.3 |
| P90/P50 | 7.6 | 7.0 | 7.4 |
| n | 4 418 | 3 698 | 4 359 |

Source: 2007 Survey of Consumer Finances.

54