

Accepted Manuscript

Notes & tips

Robust regression methods for real-time PCR

Wim Trypsteen, Jan De Neve, Kobus Bosman, Monique Nijhuis, Olivier Thas,
Linus Vandekerckhove, Ward De Spiegelaere

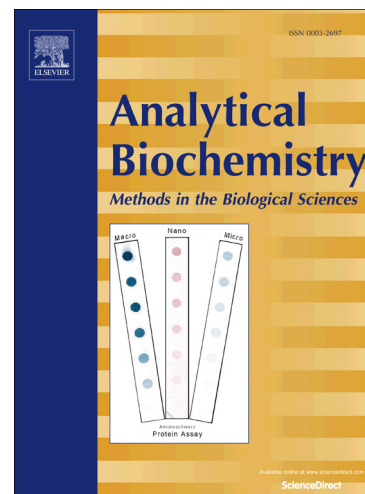
PII: S0003-2697(15)00139-6
DOI: <http://dx.doi.org/10.1016/j.ab.2015.04.001>
Reference: YABIO 12030

To appear in: *Analytical Biochemistry*

Received Date: 6 February 2015
Revised Date: 27 March 2015
Accepted Date: 2 April 2015

Please cite this article as: W. Trypsteen, J. De Neve, K. Bosman, M. Nijhuis, O. Thas, L. Vandekerckhove, W. De Spiegelaere, Robust regression methods for real-time PCR, *Analytical Biochemistry* (2015), doi: <http://dx.doi.org/10.1016/j.ab.2015.04.001>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **Robust regression methods for real-time PCR**

2 **Authors:** Trypsteen Wim¹, De Neve Jan², Bosman Kobus³, Nijhuis Monique³, Thas Olivier^{2,4},
3 Vandekerckhove Linos¹, De Spiegelaere Ward¹

4 **Affiliations:** ¹ Department of Internal Medicine, HIV Translational Research Unit, University Ghent &
5 University Hospital Ghent, Belgium. ² Department of Mathematical Modelling, Statistics and Bio-
6 informatics, University Ghent, Belgium. ³ Department of Medical Microbiology & Virology, University
7 Medical Center Utrecht, Utrecht, The Netherlands. ⁴ National Institute for Applied Statistics Research
8 Australia (NIASRA), School of Mathematics and Applied Statistics, University of Wollongong, NSW
9 2522, Australia.

10 **Corresponding author:** Vandekerckhove Linos (*); Depintelaan 185, De Pintepark Building, University
11 Hospital Ghent, 9000 Ghent, Belgium; tel: +3293323398, fax: +3293323895;
12 linos.vandekerckhove@ugent.be

13

ACCEPTED MANUSCRIPT

14 **Abstract**

15 Current real-time PCR data analysis methods implement linear least squares regression methods for
16 primer efficiency estimation based on standard curve dilution series. This method is sensitive to
17 outliers that distort the outcome and are often ignored or removed by the end-user. Here, robust
18 regression methods are shown to provide a reliable alternative, since they are less affected by
19 outliers and often result in more precise primer efficiency estimators than the linear least squares
20 method.

21 **Keywords:** Robust regression, real-time PCR, outliers, qPCR, standard curve, PCR efficiency
22 estimation

23

ACCEPTED MANUSCRIPT

24 **Manuscript**

25 The real-time quantitative polymerase chain reaction (qPCR) is a well-established technique for
26 quantification of nucleic acids [1, 2]. Quantification of relative and/or absolute quantities is generally
27 conducted by comparing the quantitative outcome to a calibration curve made by a standard dilution
28 series. In this setup, it is essential to construct a calibration curve that accurately reflects the reaction
29 efficiency of the qPCR to estimate the concentration of the unknown samples [2].

30 In current qPCR practice, the linear least squares regression is implemented to deduce the PCR
31 reaction efficiency from the calibration curve (for details see Supplementary Material 1) [3]. This
32 method works well with an optimal standard curve preparation and minimal variation between
33 technical replicates. However, in practice, minor variations are regularly observed and these tend to
34 increase in dilutions at the lower end of detection. In addition, the least squares method is sensitive
35 to outliers, especially when these are present at the extremes of the dilution series. In general,
36 variation in standard dilution series is more frequently observed in the lower extreme, when the
37 effect of the random sampling error increases. This issue often results in the manual removal of
38 some replicates by the end-user to intuitively fit a better standard curve, making data analysis
39 subjective to the interpretation of the end-user. In this light, robust regression methods for
40 estimating the slope may be preferred as they are less susceptible to outliers and provide more
41 precise estimators for a variety of error distributions [4]. This concept was recently introduced and
42 explored in the context of real-time PCR by Orenti & Marubini, showing that a robust regression
43 method, the biweight MM estimator, could offer an alternative to the linear least squares method in
44 calibration curve calculations, especially when outliers are present [5].

45 Here, we extend this concept by comparing the least squares method with three robust regression
46 procedures (the MM estimator [5,6], the robust estimator of Theil and Sen [7,8] and the robust
47 estimator of Siegel [9]), by examining the effect of the error distribution on precision in the absence
48 of outliers and by including the least-squares method after outlier removal according to the Grubbs'

49 outlier test [10]. The Grubbs test is a standardized test for detecting and removing outliers from
50 qPCR standard curves. It tests the null-hypothesis that there are no outliers in the data versus the
51 alternative that there is at least one outlier. The largest Cq value is removed when the p-value of this
52 test is less than 0.05. This process is repeated on the reduced dataset until the p-value exceeds 0.05.

53 We consider bias caused by outliers and precision as criteria to evaluate and compare the different
54 estimators. More specifically, we examine the bias of the primer efficiency estimator when there is
55 an outlier at the extreme of the dilution series. It is desirable to have a method that is not affected by
56 a single outlier when the other dilution points indicate good primer efficiency.

57 The precision is inversely proportional to the variance of the estimator. When multiple standard
58 curve dilution series of comparable quality are available for estimating the primer efficiency, then it
59 is desirable to have similar efficiency estimates across all series. This corresponds to an estimator
60 with a small standard error or, equivalently, a high precision.

61 To illustrate the performance of the different methods for standard curve calculations, a qPCR assay
62 for the quantification of HIV DNA was performed in 8 dilution series. This nested qPCR assay uses two
63 rounds of PCR amplification of the HIV gag-region. DNA from U1 cells, which contain two HIV provirus
64 integrations per cell, were used for the 8 dilution series; for details, see Supplementary Material 2
65 [11,12]. The regression lines and corresponding efficiencies for each replicate are estimated
66 according to the different regression methods (Supplementary Material 3 & 4). All calculations were
67 performed using the statistical software environment R [13].

68 A first distinction between linear and robust regression methods can be made by examining a
69 dilution series with and without an outlier (Figure 1). In case no outliers are present (Figure 1A), all
70 four methods provide a comparable and accurate estimate of the efficiency. However, when an
71 outlier is present (Figure 1B), the least squares method results in a decreased primer efficiency
72 estimate, while the estimated efficiencies of the robust methods remain largely unaffected by the
73 outlier.

74 To examine this in more detail, data simulations were performed to mimic the experimental setups
75 with and without outliers. Seven distinct calibration curves are considered and replicated twice.
76 Quantification cycles are simulated according to a linear model with intercept 26.5 and slope -3.553
77 corresponding to a primer efficiency of 91.2% (see the Supplementary Material 1 for more details on
78 the linear model and the primer efficiency). To simulate the errors, we consider the normal
79 distribution with mean zero and standard deviation 0.19 which corresponds to the estimated
80 standard deviation of the error distribution of the data in Figure 1A.

81 Table 1A shows the average estimated efficiency based on 10000 Monte-Carlo simulations. Without
82 the outlier, all estimators are unbiased. However, when a single outlier is present, the least squares
83 estimator is biased and, on average, underestimates the efficiency by 21%. On the other hand, the
84 MM estimator remains unbiased while the Theil-Sen and Siegel estimators underestimate the
85 efficiency by approximately 1%.

86 A second distinction can be made by examining the variance of the efficiency estimators. For the real
87 data, the robust regressions have lower standard deviations for the estimated efficiency compared to
88 the linear least squares regressions (Supplementary Material 4). These standard deviations may
89 suggest that the least squares estimators are less precise as compared to the robust estimators so
90 that robust regression will likely estimate the true efficiency more accurately. This is of great
91 importance while newly designed qPCR primer pairs are often only tested once or twice on a
92 standard dilution series of reference material.

93 To examine whether the presence of outliers are solely responsible for this observation, data
94 simulations were performed as described above, but only to mimic the experimental setup without
95 outliers so that all four estimators are unbiased. In addition, the performances of the estimators are
96 tested over different error distributions, as the underlying distribution of the error is unknown. The
97 following error distributions were included: normal, student t with 3 degrees of freedom, lognormal,
98 Gumbel and Laplace. They were all standardized to mean zero and standard deviation 0.19 to make

99 the results comparable (Table 1B). These choices include symmetric, heavy tailed, and skewed
100 distributions.

101 For the five error distributions considered, the MM and Theil-Sen estimator outperform the least
102 squares estimators except for the normal distribution for which the least squares estimator has a
103 slightly better performance. Hence, even without outliers, robust regression methods can produce
104 more precise estimates.

105 For the normal distribution, the least squares estimator is most precise. This can be expected since
106 the least squares estimator corresponds to the maximum likelihood estimator for normal distributed
107 errors. In comparison with the least squares estimator, the MM, Theil-Sen and Siegel estimators
108 result in an increased standard error of 5%, 5% and 17% respectively. For the other error
109 distributions, however, the MM and Theil-Sen estimators have an increased precision over the least
110 squares estimator. For the lognormal distribution, for example, the standard error of the least
111 squares estimator is more than twice the standard error of the MM and Theil-Sen estimators. The
112 Siegel estimator has a superior precision over the least squares for the t-distribution, lognormal and
113 Laplace distribution, but underperforms as compared to the MM and the Theil-Sen estimators.

114 In summary, we have demonstrated that robust regression estimators are less affected by outliers
115 and often prove to be more precise for estimating PCR efficiency compared to standard linear
116 regression with least squares estimation. Especially the MM and Theil-Sen estimators seem to be
117 appropriate for primer efficiency estimation. Therefore, the implementation of robust regression
118 methods in qPCR analysis would provide a reliable alternative. Single outlying dilutions can
119 effectively introduce bias to efficiency estimates, particularly when these occur in the extremes of
120 the dilution series where stochastic sampling effects at the lower levels of detection may increase
121 the variation. However, we do want to stress that outlying data points can be an important indication
122 of poor pipetting and low quality data of the entire PCR run. Therefore, robust regression methods
123 should not be used to falsely expand the dynamic range of the linear interval of the PCR assay. The

124 Cq values at which outlying data points are observed, should still be considered outside the dynamic
125 range of the assay in the given PCR run.

126 Acknowledgements

127 Financial support was provided by the King Baudouin Foundation (Grant 2010-R20640-003) and the
128 HIV-ERA Eranet project EURECA: SBO-IWT (Grant 130442). Prof. Linos Vandekerckhove is supported
129 by the Research Foundation – Flanders (FWO; Grant 1.8.020.09.N.00). Kobus Bosman is supported by
130 the AIDS FONDS (project 2013034), Monique Nijhuis is supported by the amfAR Research Consortium
131 on HIV Eradication (ARCHE) (EPISTEM, project (108930) and the Netherlands Organization for
132 Scientific Research VIDI (grant number 91796349). Jan De Neve and Olivier Thas gratefully
133 acknowledge the IAP research network grant no. P7/06 of the Belgian government (Belgian Science
134 Policy). The research fits within the N2N Multidisciplinary Research Partnership of Ghent University
135 (01MR0310W).

136 References

- 137 [1] C. A. Heid, J. Stevens, K. J. Livak, P. M. Williams, Real time quantitative pcr, *Genome Res* 6
138 (1996) 986-994.
- 139 [2] S. A. Bustin, V. Benes, J. A. Garson, J. Hellemans, J. Huggett, M. Kubista, R. Mueller, T. Nolan,
140 M. W. Pfaffl, G. L. Shipley, J. Vandesompele, C. T. Wittwer, The miqe guidelines: Minimum
141 information for publication of quantitative real-time pcr experiments, *Clin Chem* 55 (2009)
142 611-622.
- 143 [3] J. M. Ruijter, M. W. Pfaffl, S. Zhao, A. N. Spiess, G. Boggy, J. Blom, R. G. Rutledge, D. Sisti, A.
144 Lievens, K. De Preter, S. Derveaux, J. Hellemans, J. Vandesompele, Evaluation of qpcr curve
145 analysis methods for reliable biomarker discovery: Bias, resolution, precision, and
146 implications, *Methods* 59 (2013) 32-46.
- 147 [4] P.J. Rousseeuw, A.M. Leroy, *Robust regression and outlier detection*, vol 589, John Wiley and
148 Sons, 205.
- 149 [5] A. Orenti, E. Marubini, Performance of robust regression methods in real-time polymerase
150 chain reaction calibration, *The International journal of biological markers* 29 (2014) e317-
151 327.
- 152 [6] V.J. Yohai, High Breakdown-point and High Efficiency Estimates for Regression, *The Annals of*
153 *Statistics*, 15 (1987), 642-65.
- 154 [7] H. Theil, , A Rank-Invariant method of linear and polynomial regression analysis, I, II, and III,
155 *Nederlandsche Akad. van Wetenschappen Proc.* 53 (1950) 386-92, 521-5, 1397-412.

- 156 [8] P.K. Sen, Estimates of the Regression Coefficient Based on Kendall's Tau, Journal of the
157 American Statistical Association 63 (1968) 1379-1389.
- 158 [9] A.F. Siegel, Robust Regression Using Repeated Medians, Biometrika (1982) 242-244.
- 159 [10] M.J. Burns, G.J. Nixon, C.A. Foy, N. Harris, Standardisation of data from real-time quantitative
160 PCR methods – evaluation of outliers and comparison of calibration curves, BMC
161 Biotechnology 7 (2005) 5-31.
- 162 [11] W. Hu, R. Kaminski, F. Yang, Y. Zhang, L. Cosentino, F. Li, B. Luo, D. Alvarez-Carbonell, Y.
163 Garcia-Mesa, J. Karn, X. Mo, K. Khalili, RNA-directed gene editing specifically eradicates latent
164 and prevents new HIV-1 infection, Proceedings of the National Academy of Sciences of the
165 United States of America 111 (2014) 11461-11466.
- 166 [12] T.M. Folks, J. Justement, A. Kinter, S. Schnittman, L. Orenstein, G. Poli, A.S. Fauci,
167 Characterization of a promonocyte clone chronically infected with HIV and inducible by 13-
168 phorbol-12-myristate acetate, Journal of Immunology 140 (1988) 1117-1122.
- 169 [13] R Core Team, R: A language and environment for statistical computing, R Foundation for
170 Statistical Computing, Vienna, Austria, (2014) URL <http://www.R-project.org/>.

171 **Figure 1: Estimated regression lines and estimated efficiencies for two of the eight replicates.**
172 Comparison of the five regression estimates without an outlying dilution (A) and in case of an
173 outlying dilution at the lower extreme of the standard curve (B). The least squares (black line, --), the
174 least squares after removal of outliers with the Grubbs test (red line, ...), robust MM (green line, ·-·),
175 the robust Theil-Sen (blue line, —) and the robust Siegel (turquoise line, - - -) methods show
176 equal efficiency in the upper panel (A). The least squares estimate is heavily affected by the outlier in
177 the lower panel (B), while the estimated efficiencies of the robust methods and the least squares
178 after performing Grubbs' test for removing the outlier remain comparable across the replicated
179 experiments.

180

ACCEPTED MANUSCRIPT

181 **Table 1: Results of the simulation study in which the true efficiency equals 91.2%.** A: averages of
 182 the estimated efficiencies (in %) for the five estimators in the presence and absence of an outlier. B:
 183 standard deviations of the estimated efficiencies for the five estimators applied to the simulated data
 184 without outlier (for which all estimators are unbiased). Smaller values indicate more precise
 185 estimators. All results are obtained based on 10000 Monte-Carlo simulations.

186

A	LS	LS + Grubbs	MM	Theil-Sen	Siegel
Without outlier	91.2	91.2	91.2	91.2	91.2
With outlier	70.0	91.2	91.2	90.1	90.5

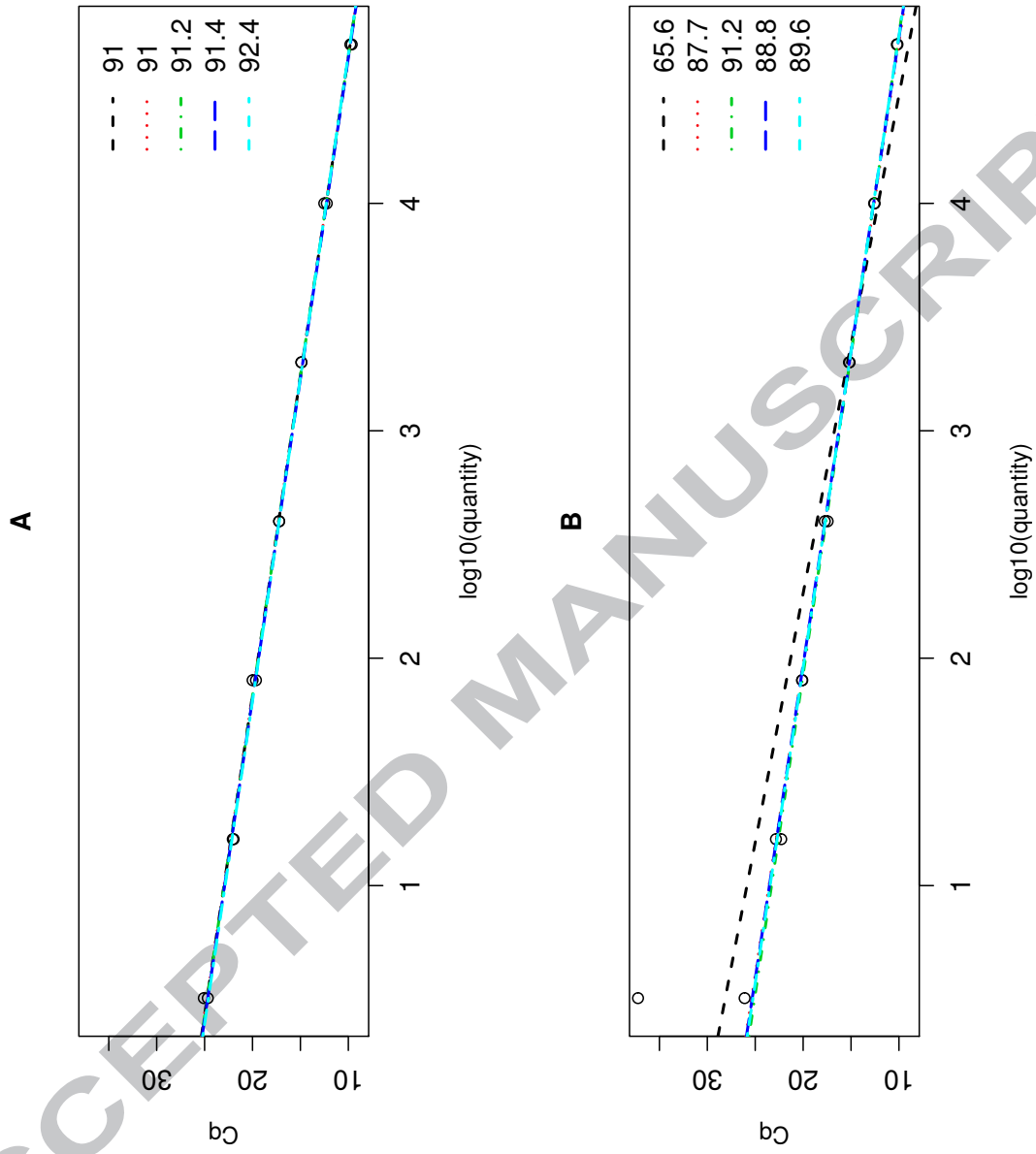
187

B	LS	LS + Grubbs	MM	Theil-Sen	Siegel
Normal	1.27	1.27	1.33	1.34	1.48
Student t	1.26	1.28	0.99	1.00	1.05
Lognormal	1.27	1.27	0.64	0.62	0.68
Gumbel	1.27	1.27	1.22	1.21	1.34
Laplace	1.27	1.27	1.17	1.16	1.20

188

189

190



191