

Volume 116, Number 6, June 2016

British Journal of Anaesthesia **116** (6): 733–736 (2016) doi:10.1093/bja/aew110

EDITORIALS

Academic assessment of arterial pulse contour analysis: missing the forest for the trees?

A. Reisner*

Harvard Medical School, Department of Emergency Medicine, Massachusetts General Hospital, Zero Emerson Place, Suite 3B, Boston, MA 02114, USA

*E-mail: areisner@partners.org

In this issue of the British Journal of Anaesthesiology, Montenij and colleagues¹ provide a thoughtful review of analytic methods for comparing cardiac output measurement methods, with focus on arterial pulse contour analysis methods that are intended to measure cardiac output. This review is a welcome addition to the literature, as such comparative investigations are common-place, and often without optimal rigor.

At the same time, I hold a concern about another deficiency in pulse contour investigations. Despite 414 'pulse contour' AND 'cardiac output' articles currently indexed by PubMed, a century of academic reports describing various pulse contour methods² and decades of commercial sales, it remains uncertain whether pulse contour methods provide more sensitive and specific indicators about circulatory decompensation than routine use of blood pressure (bp) and heart rate monitoring – let alone whether such technology leads to improved patient outcomes. It can be argued that conducting future studies that continue to merely compare one cardiac output measurement technique against another, risks missing the forest for the trees.

To be sure, management of the tenuous patient – stable but with minimal physiological reserve, and with a high risk of decompensation, as might occur after haematemesis from esophageal varices, or during an invasive procedure – is a challenge. The vigilant clinician monitors tenuous patients carefully, to respond to any deterioration while avoiding unnecessary and excessive intervention. A conundrum may occur when the arterial bp drifts down, which might indicate deterioration, such as new, dangerous blood losses. Yet as often as not, this is relatively benign, a lessening of vasoconstriction as the patient becomes more relaxed as a result of medication or time. The vigilant clinician must distinguish between these two very different physiological circumstances, and the stakes are high. Accordingly, there has been interest in techniques for non-invasive monitoring cardiac output, such as pulse contour analysis, as cardiac output is a cardinal metric of circulatory adequacy and, as the dividend of the bp-to-central venous pressure gradient, also yields total peripheral resistance (TPR), a measure of vasoconstriction.³

If pulse contour analysis for measuring cardiac output is truly reliable, then it should be uniformly used as the standard-of-care for tenuous patients. If the technique is inaccurate, then it is only an illusion that the patient's cardiac output and vascular tone are being carefully monitored, unreliable information that gives a false – and possibly dangerous – sense of security when managing tenuous patients.

The rationale behind pulse contour analysis

For a masterly treatment of the principles underlying pulse contour analysis, one may consult the textbook 'McDonald's Blood Flow in Arteries'.⁴ The vast majority of pulse contour methods estimate volumetric flow in the aortic root, which equals cardiac output. As a matter of basic physics, note that it is <u>not</u> the pressure <u>wave</u> but the pressure <u>aradient</u> that impels fluid to flow in blood vessels. In other words, it is the <u>difference</u> of pressures in a segment of artery, <u>upstream</u> vs downstream, that accelerates/ decelerates the <u>pulsatile blood</u> within that segment. It is simply

© The Author 2016. Published by Oxford University Press on behalf of the British Journal of Anaesthesia. All rights reserved. For Permissions, please email: journals.permissions@oup.com

impossible to compute flow with only one pressure wave: computing the gradient cannot be done precisely without a second pressure measurement. Most pulse contour methods address this conundrum by relying on <u>probabilistic</u> relationships between the upstream pressure wave and the downstream pressure wave. For instance, one method of calculating flow would be to assume that the downstream pressure waveform is similar to the upstream pressure, aside from a small time delay. From this assumption, pressure gradients can be estimated and flow computed.

The first challenge for pulse contour analysis is that the downstream waveform is not, in fact, the same shape as the upstream waveform, because the entire pressure waveform is not wholly moving downstream. Instead, there is a primary wave moving downstream, while there are smaller pressure waves that move upstream: reflected pressure waves from distal vascular junctures that travel in the retrograde direction. These reflected waves increase the amplitude of an arterial waveform, but they actually create retrograde pressure gradients that decelerate the blood and retard flow. Larger bp waveforms do not always correspond to greater forward flow! This is one fundamental challenge to pulse contour analysis, and different pulse contour methods use different techniques, usually statistical corrections based on either patient characteristics or some property of the shape of the bp waveform, to try to circumvent this complication.

There is a second challenge. While the gradient of pressure determines the magnitude of acceleration/deceleration, it is the <u>diameter</u> of the <u>vessel</u> that dictates the actual <u>volume</u> of <u>blood</u>. Pulse contour methods must use some technique to address this complication, such as a calibration of volume against another reference, or relying on <u>probabilistic relationships</u> between age and gender and the likely size and pulsatile <u>compliance</u> of the arterial <u>vessel</u>.

There is also a third major analytic challenge. It is only within the aortic <u>root</u> that <u>flow</u> equals <u>cardiac</u> <u>output</u>, whereas the arterial waveform is usually <u>measured</u> somewhere in the <u>periphery</u>. Pulse contour methods must use some technique to <u>estimate</u> flow in the proximal aorta using a pressure waveform measured in the periphery. Again, a common approach is to use <u>probabilis-</u> tic relationships between those waveforms. (One approach is to use a generalized transfer function, which is a mathematical manipulation based entirely on <u>probabilistic</u> relationships between peripheral and central waveforms⁵).

The crux of the matter

There is indeed a physical causal relationship between the arterial pressure waveform and cardiac output. However, taking the three analytic challenges together, it is also clear than quantifying cardiac output from pulse contour analysis must rely on probabilistic relationships (e.g. the likely relationship between the central and peripheral arterial waveform; the likely relationship between the patient's age, gender, etc. and the size and compliance of the patient's aorta; or the likely relationship between the upstream and downstream pressure waves that determine the flow-determining pressure gradient). These 'likely relationships' are purely probabilistic; they are observed in the majority of cases, but not all cases. Conceptually, it is no different from relying on a patient's weight to estimate her height: a reasonable estimate can be made for many individuals, but there is likely a subset for whom the relationship will be invalid. Which means that the estimated cardiac output by pulse contour may be accurate, except when it isn't.

This doesn't invalidate pulse contour analysis. We clinicians are accustomed to relying on probabilistic relationships when we assess our patients' haemodynamics. When mean arterial pressure (MAP) is falling, we know that it typically represents failing circulation. Or when the patient has a large pulse pressure, we assume that the patient probably has a large stroke volume. Yet sometimes these probabilistic relationships are invalid (e.g. patients with low MAP who are not in shock but are merely vasodilated). It is because our routine measures, such as MAP and pulse pressure, can be clinically ambiguous that we seek superior, less ambiguous non-invasive measures.

Returning to pulse contour analysis: the motivation to incorporate this technology into our practice is because we know that routine bp is not always reliable in assessing circulatory state. Yet pulse contour cardiac output also relies on a set of probabilistic relationships that may be invalid for some subset of clinical situations. Is pulse contour analysis superior to routine vital signs monitoring? Or does it provide false reassurance by continuously displaying a cardiac output estimate that is not always reliable? Is it indeed superior to routine monitoring? In my opinion, after more than several decades, this question is not answered.

Open questions within the literature

The academic literature regarding pulse contour analysis is dominated by method comparison studies (i.e. comparing cardiac output from pulse contour analysis vs a reference method). Method comparison studies do not answer the following question: 'should I use technology X for patient Y.' The emphasis on '95% confidence intervals' can mask serious problems that can occur under certain circumstances: pulse contour method, overly reliant on probabilistic relationships, might yield wildly and systematically inaccurate cardiac output in a subset of patients with atypical physical or physiological properties. Focus on the majority of cases who, by definition, fall within the 95% confidence interval, and treating errors as if they are just 'random', means that major failures of these techniques are treated as nothing more than 'outliers' (i.e. unpredictable statistical flukes).

Yet it is very possible that pulse contour analysis fails in a predictable way under predictable conditions (and those conditions may or may not be commonplace in any given published study). It would be valuable to determine if there are patients in whom pulse contour analysis predictably fails so that we may learn the most about the technologies' true capabilities and pitfalls. Consider pulse oximetry, by analogy: we know not to rely on pulse oximetry if there are haemoglobinopathies or after methylene blue, and we know it is less reliable given poor skin perfusion or bright ambient lights. We must focus on defining any non-random sources of error for each and every investigational cardiac output method that we hope to use on patients.

Method comparison studies¹ are only a rudimentary way of assessing pulse contour analysis. Other essential questions include, how frequently does clinical management change when guided by pulse contour analysis rather than routine monitoring, and are overall prospective outcomes improved? Or, can pulse contour analysis predict the patient's future physiological state better than routine methods involving bp and heart rate alone? (A useful schema for diagnostic test assessment includes technical efficacy, diagnostic accuracy efficacy, diagnostic thinking efficacy, therapeutic efficacy, clinical outcome efficacy, and societal efficacy⁶). There are a relatively small number of published prospective outcomes trials involving pulse contour analysis for cardiac output monitoring. Of those outcomes trials, many are suboptimal, comparing management using pulse contour analysis vs *ad* hoc care. Ad hoc care, in which there are no explicit expectations in how the control group is managed, is often inferior to *any* rigorous protocol. Consider that the <u>bispectral</u> index monitor (for anaesthesia monitoring⁷) and the continuous fiber optic central venous oximetry monitor (for sepsis resuscitation⁸) were both associated with significantly superior outcomes than *ad* hoc care. However, when those technologies were later <u>compared</u> against <u>alternative monitoring</u> methodologies, <u>neither</u> novel technology was found to be <u>superior</u>.^{9 10}

A second problem with many outcomes trials involving pulse contour analysis is that they use a *bundle* of technologies from the vendor, such as **pulse contour cardiac output** and stroke volume variation metrics for predicting volume responsiveness. Certainly, if using a bundle of technologies can be shown to improve patient outcomes, that finding is noteworthy and may be practice-changing. However, study designs involving bundles do not reveal which technologies within the bundle are reliable and which are not.

Overall, there exist a rather limited number of studies investigating whether or not patients experience superior outcomes using pulse contour analysis monitoring of cardiac output, and not nearly enough is known about a technology that has been sold and used in patient care for decades. Searching the Cochrane Library, there is only one meta-analysis involving pulse contour analysis (which concluded 'an absence of evidence that fluid optimization strategies improve outcomes for participants undergoing surgery for [proximal femur fracture' and 'length of hospital stay may be improved, but lack of good quality data leaves uncertainty').¹¹ Additional outcomes investigations, including replication of successful pilot studies without industrysponsorship, should be encouraged for anyone interested in the academic assessment of pulse contour methods.

Closing remarks

The current landscape is a generally poor understanding of pulse contour analysis products' clinical value. Searching the Pulsion. com website,¹² I cannot find any detailed explanation for their pulse contour analysis algorithm that addresses the analytic challenges that were discussed above. The FloTrac website¹³ is similarly vague, merely stating broadly that: 'Cardiac output is correlated with the variance between systolic and diastolic pressure. Real-time analysis of waveform characteristics is also integrated, compensating for changes in vascular physiology affecting the pressure waveform.' Retia Medical offers a different approach to pulse contour analysis, seeking to estimate the rate at which blood drains from the arterial tree by looking at the arterial pressure over long time intervals (rather than estimating volume of each systole)¹⁴ but their website does not provide substantial detail, either.¹⁵ We clinicians would never treat patients with pharmacologic agents with active ingredients so poorly described and understood.

There is insufficient published evidence that pulse contour methods, which offer cardiac output estimates based on a set of probabilistic assumptions, are safe, or that they offer significant advantage over careful monitoring using routine haemodynamic parameters. Nor is it clear which commercially sold methods are superior to the others. There is generally weak understanding about the specific conditions under which each method is reliable vs erroneous. The fact that pulse contour analysis generally correlates well with cardiac output references measurements (as demonstrated by many cardiac output method comparison studies) is **not**, to me, adequately **encouraging**, because **changes in MAP also generally correlate well with changes in cardiac output**.¹⁶ Therefore, correlation with cardiac output is not in-and-of-itself evidence that pulse contour analysis is better than routine use of bp and heart rate monitoring.

The promise of pulse contour analysis is clear: there is a wealth of information within the waveform, and it seems reasonable that rigorous analysis should yield superior information for decision-making than rudimentary metrics such as systolic bp and MAP. In the near future, one hopes that these open questions can be addressed with trials that unambiguously affirm the clinical value of pulse contour analysis for patient management, while clinicians will use the technologies with a deep understanding of its underlying principles and the evidence that it is useful for a given clinical application.

Declaration of interest

A.T.R. has received funding from Nihin-Kohden for decision support in sepsis management, and holds patents related to electronic decision support.

References

- 1. Montenij LJ, Buhre W, Jansen JRC, Kruitwagen C, De Waal E. Methodology of method comparison studies evaluating the validity of cardiac output monitors: a stepwise approach and checklist. Br J of Anaesth 2016; **116**: 750–8
- Erlanger J, Hooker DR. An experimental study of blood pressure and of pulse-pressure in man. Johns Hopkins Hosp Rep 1904; 12: 145–378
- Soble JS. In search of the Holy Grail. Critical Care Medicine 1998; 26: 1953–4
- O'Rourke MF, Nichols W, Vlachpoulos C. McDonald's Blood Flow in Arteries 6th Edition: Theoretical, Experimental and Clinical Principles. London: Hodder Arnold, 2011
- Gallagher D, Adji A, O'Rourke MF. Validation of the transfer function technique for generating central from peripheral upper limb pressure waveform. Am J Hypertens 2004; 17(11 Pt 1): 1059–67
- Pearl WS. A Hierarchical Outcomes Approach to Test Assessment. Annals of Emergency Medicine 1999; 33: 77–84
- Myles PS, Leslie K, McNeil J, Forbes A, Chan MT. Bispectral index monitoring to prevent awareness during anaesthesia: the B-Aware randomised controlled trial. *Lancet* 2004; 363: 1757–63
- Rivers M, Nguyen B, Havstad S, et al. Early goal-directed therapy in the treatment of severe sepsis and septic shock. N EnglJ Med 2001; 345: 1368–77
- Yealy DM, Kellum JA, Huang DT, et al. A randomized trial of protocol-based care for early septic shock. N Engl J Med 2014; 370: 1683–93
- Avidan M, Zhang L, Burnside B, Finkel K, Searleman A, Selvidge J. Anesthesia awareness and the bispectral index. N Engl J Med 2008; 358: 1097–108
- Brammar A, Nicholson A, Trivella M, Smith A. Perioperative fluid volume optimization following proximal femur fracture. Cochrane Database of Syst Rev 2013; 9: CD003004
- Pulsion Medical Systems. Available from http://www.pulsion. com/international-english/home/ (accessed 29 January 2016)

- Edwards Lifesciences. Available from http://www.edwards. com/products/mininvasive/Pages/flotracfaqs.aspx (accessed 29 January 2016)
- Mukkamala R, Reisner AT, Hojman HM, Mark RG, Cohen RJ. Continuous cardiac output monitoring by peripheral blood pressure waveform analysis. *IEEE Transactions on Biomedical Engineering* 2006; 53: 459–67
- Retia Medical. Product page for hemodynamic monitor. Available from http://www.retiamedical.com/products/ hemodynamic-monitor (accessed 29 January 2016)
- Sun JX, Reisner AT, Saeed M, Heldt T, Mark RG. The cardiac output from blood pressure algorithms trial. Critical Care Medicine 2009; 37: 72–80

British Journal of Anaesthesia **116** (6): 736–738 (2016) doi:10.1093/bja/aew149

Applied cardiovascular physiology in theatre: measuring the cardiovascular effects of propofol anaesthesia

M. R. Pinsky^{1,2,*}

¹ Department of Critical Care Medicine and Anaesthesiology, University of Pittsburgh, Pittsburgh, PA, USA, and
² Department of Anaesthesiology, University of California, San Diego, CA, USA

*E-mail: mclaughlinbr@upmc.edu, pinskymr@upmc.edu

Cardiovascular homeostasis is a complex and beautiful interplay between the functional differences between various vascular circuits in the body and their tissue's metabolic demand, the physical nature of the endothelial barrier to fluid flux, the circulating blood volume, and reflex-mediated autonomic tone. When at rest, as occurs during anaesthesia, basal metabolic demand is both constant and low. Thus, impairments in autoregulation or sudden decreases in blood volume, as may happen during surgery, are thankfully less detrimental to tissue wellness than might otherwise be the case under conditions of metabolic stress. However, such physiologic reserve though comforting to the anaesthetist and forgiving to the patient, has clearly defined limits. Anaesthesia by its nature decreases central nervous system activity and by default, impairs autonomic responsiveness and at high enough concentrations impairs vascular tone and cardiac contractility. These concepts form the basis for anaesthetic selection in specific patient groups. But mostly all these considerations have focused on the left ventricle (LV) and arterial tone, ignoring venous return by simply placating it with increased fluid resuscitation, vasopressor infusion and/or decreased concentration of anaesthesia if the patient becomes hypodynamic.

However, the circulation is much more interactive in its components defining cardiac output than those described by left ventricular preload and contractility and arterial pressure and arterial vasomotor tone. Fundamental principles of cardiovascular physiology, as originally described by Guyton and colleagues¹ more than 50 yr ago,¹ identified venous return as the primary determinant of cardiac output and that LV function is remarkably insensitive in defining this level of flow, only the required backpressures needed for that flow. We collectively argued these points relative to cardiopulmonary bypass surgery in a physiologic commentary.² Until recently, just knowing that venous return was the primary determinant of cardiac output did little to help the bedside clinician manage complex and changing surgical patients. One understood that mean circulatory filling pressure (Pmcf) was the best surrogate for effective circulating blood volume, but its measure and its own determinants were difficult to ascertain at the bedside and nearly impossible to measure repeatedly over time. The effective circulating blood volume represents a balancing act between total circulating blood volume, blood flow distribution amongst various organs with varying degrees of capacitance and unstressed volume, and the resistance to venous return (RVR), which has more of a conductance determinant to its value that actual physical resistive.3 Importantly, multiple lines of investigation have led to the development of several methods to quantify Pmcf at the bedside using only arterial pressure, central venous pressure (CVP), and cardiac output. A detailed review of these various techniques is found elsewhere.⁴ However, presently three techniques are readily available and can be used for the bedside assessment of venous return.

The first approach uses an analogue estimate of Pmcf by assuming a constant proportion of compliance and resistances within the arterial and venous circuit.⁵ We recently validated this breath-by-breath analogue approach in a canine model during normal and endotoxic shock state.⁶ Using this analogue approach Cecconi and colleagues⁷ examined the effect of fluid boluses on Pmcf, the driving pressure for venous return (Pmcf-CVP), and cardiac output in a large postoperative surgical patient population. They showed that fluid loading universally increased Pmcf, if only transiently, and unaltered RVR. However, for cardiac output to increase the driving pressure for venous return also needed to increase. Thus, if fluid loading did not increase cardiac output, CVP increased, whereas in those whose cardiac output increased CVP remained stable. The observation that volume loading does not alter RVR has been known for more than 30 yr,⁸ and is the basis for increases in CVP during fluid loading being a 'stopping rule' for fluid infusion therapy.⁹

doi: 10.1093/bja/aew094 Review Articles

REVIEW ARTICLES

Methodology of method comparison studies evaluating the validity of cardiac output monitors: a stepwise approach and checklist[†]

L. J. Montenij¹, W. F. Buhre², J. R. Jansen³, C. L. Kruitwagen⁴ and E. E. de Waal^{1,*}

¹Department of Anaesthesiology, Intensive Care and Emergency Medicine, University Medical Centre Utrecht, Heidelberglaan 100, 3584 CX Utrecht, The Netherlands, ²Department of Anaesthesiology and Pain therapy, Maastricht University Medical Centre, P. Debeyelaan 25, 6229 HX, Maastricht, The Netherlands, ³Department of Intensive Care, Leiden University Medical Centre, Albinusdreef 2, 2333 ZA, Leiden, The Netherlands, and ⁴Department of Biostatistics and Research Support, Julius Centre, Universiteitsweg 100, 3584 CG Utrecht, The Netherlands

*Corresponding author. E-mail: e.e.c.dewaal@umcutrecht.nl

Abstract

The validity of each new cardiac output (CO) monitor should be established before implementation in clinical practice. For this purpose, method comparison studies investigate the accuracy and precision against a reference technique. With the emergence of continuous CO monitors, the ability to detect changes in CO, in addition to its absolute value, has gained interest. Therefore, method comparison studies increasingly include assessment of trending ability in the data analysis. A number of methodological challenges arise in method comparison research with respect to the application of Bland–Altman and trending analysis. Failure to face these methodological challenges will lead to misinterpretation and erroneous conclusions. We therefore review the basic principles and pitfalls of Bland–Altman analysis in method comparison studies concerning new CO monitors. In addition, the concept of clinical concordance is introduced to evaluate trending ability from a clinical perspective. The primary scope of this review is to provide a complete overview of the pitfalls in CO method comparison research, whereas other publications focused on a single aspect of the study design or data analysis. This leads to a stepwise approach and checklist for a complete data analysis and data representation.

Key words: cardiac output; trends; validation studies

Method comparison research aims to evaluate the validity of a new monitor against an established reference technique, and is of specific interest in cardiac output (CO) monitoring.¹² After establishing validity, other types of research are needed to evaluate the extent to which new monitors alter haemodynamic management, effects on patient outcome, and cost-effectiveness. Method comparison studies face a number of methodological challenges. A number of reviews have been published, most of them discussing a component of the application of Bland-Altman analysis in this setting.^{3–7} Despite these reviews, many

studies do not meet a number of fundamental principles.³ ⁸ This may lead to incorrect conclusions and erroneous applications in clinical practice. This review therefore aims to provide a complete overview of the methodological considerations in method comparison studies concerning new CO monitors. Each component of the study design, data analysis, or data interpretation is followed by a recommendation. In addition, we focus on evaluation of trending ability, which has become increasingly important with the emergence of continuous systems.^{9 10} Current methods to analyse trending ability have a number of

© The Author 2016. Published by Oxford University Press on behalf of the British Journal of Anaesthesia. All rights reserved. For Permissions, please email: journals.permissions@oup.com

[†] This Article is accompanied by Editorial Aew110.

limitations.¹¹ As an alternative, the concept of 'clinical concordance' and a corresponding error grid method for evaluation of trending ability from a clinical perspective is introduced. Finally, the methodological issues are summarized, resulting in a stepwise approach and checklist for CO method comparison research. The use of this checklist could lead to a more complete and homogeneous presentation of data, which may facilitate systematic reviews and meta-analyses in the future.

Bland-Altman analysis: concept

Each new CO monitor should be evaluated for its accuracy and precision; accuracy refers to the ability to measure CO close to

its true value, whereas precision refers to the spread of repeated measurements (Fig. 1A). Measurement of the 'true' CO is extremely difficult in clinical practice, and reference techniques can provide only an approximation.^{1 2} This problem can be handled in part using Bland–Altman analysis.^{12–14} This method evaluates agreement between two measurement techniques, rather than validating the experimental technique against a perfect reference. As a result, only conclusions about interchangeability between the experimental and reference technique can be drawn. Bland–Altman analysis determines the bias, or mean difference between the experimental and reference technique, as a measure of accuracy.^{12–14} As a measure of precision, the 95% limits of agreement (LOA) are used (Fig. 1B). The LOA are generally



Fig 1 Accuracy and precision and the relation with bias and the limits of agreement as determined with Bland–Altman analysis. (A) Accurate measurements are close to the true value, irrespective of the spread of the measurements; in contrast, precise measurements are close to each other, irrespective of their deviation from the true value. Valid cardiac output monitors are both accurate and precise. (B) In Bland–Altman plots, accurate cardiac output monitors show a bias (continuous line) close to the line 'x=0', whereas precise monitors show limits of agreement close to the bias (dotted lines). CO_{exp}, cardiac output of the experimental technique; CO_{ref}, reference cardiac output; LOA⁺, upper limit of agreement; LOA⁻, lower limit of agreement.

determined as:

$$LOA = (bias) \pm t_{\alpha,n-1} * (sD)$$

in which SD is the standard deviation of the differences, *n* the sample size, and $t_{\alpha,n-1}$ the t-value corresponding to the degrees of freedom (n-1) and a type I error (α) of 0.05. The LOA therefore represent the limits enclosing 95% of the differences. The bias and LOA can be depicted in a Bland–Altman plot (Fig. 1B). The mean error or percentage error is calculated as follows:

Mean error (%) =
$$100\% * t_{\alpha,n-1} * \frac{(sD)}{(mean CO)}$$

Consequently, the mean error is a measure of interchangeability relative to the underlying CO and therefore a more appropriate parameter to compare the results from different studies. For calculation of the LOA and mean error, a t-value of 1.96 is often used. Strictly speaking, this value holds true only in infinitely large sample sizes. It is advisable to use correct t-values in small studies (e.g. <20 subjects), as a value of 1.96 will underestimate the real LOA and mean error.

Pitfalls in the application of Bland-Altman analysis

Bland–Altman analysis has a number of important pitfalls, which are discussed in the next sections, each followed by a recommendation. These recommendations are summarized in a checklist (Table 1).

Normal distribution

The differences between the experimental and reference technique should be normally distributed. Usually, this will be the case, even if the individual CO measurements with the experimental or reference technique do not follow a normal distribution.¹⁴ If not, a straightforward non-parametric approach is available.^{14 15} Normal quantile–quantile (QQ) plots or histograms of the differences provide a visual check of normality.¹⁶ In addition, the Kolmogorov–Smirnov or Shapiro–Wilk test can be applied. Nonetheless, small studies may pass these tests because of insufficient statistical power to demonstrate non-normality. In contrast, large studies tend to be tested non-normal even if the deviation from a normal distribution is small.¹⁶

Recommendation

Check the differences between the experimental and reference technique for normality by combining a visual check and statistical test.

Proportionality

The bias and LOA are meaningful estimates only if they are uniform over the range of measurements.¹⁴ If the difference between the techniques increases with an increase in CO, the bias will be overestimated in the low-CO range and underestimated in the high-CO range. This effect is called proportional bias and can be quantified by plotting a regression line in the Bland–Altman plot. If the slope of this line differs significantly from zero, proportional bias is present.^{14 17} Nonetheless, in small studies, proportional bias cannot be ruled out because these studies may lack the statistical power to demonstrate this significant difference. Regression analysis should therefore be accompanied

Table 1 A stepwise approach and checklist to the design, data analysis, and data interpretation of cardiac output method comparison studies. CI, confidence interval; CO, cardiac output; LOA, limits of agreement; 4Q, four-quadrant; TDCO, thermodilution cardiac output; TPCO, transpulmonary thermodilution cardiac output

Study phase	Торіс	Checklist item
Design	Measurement protocol	Create a protocol for the timing and recording of CO measurements, considering haemodynamic fluctuations, dependence of paired measurements, and the response time of (continuous) systems
	Criteria for agreement	Define criteria for acceptable bias and LOA or mean error, depending on the clinical context
	Sample size	Consider a sample size calculation (suggested method in Supplementary Appendix B or method by Bland), ²¹ or assess the appropriate sample size based on historical data
	Reference technique	Choose a highly precise reference technique (e.g. TDCO or TPCO)
Data analysis	Normal distribution	Check whether the differences are normally distributed by combining a visual check and statistical test
	Bland–Altman analysis	Calculate the bias, LOA, mean error, and their corresponding 95% CIs, using correct t-values A correction for the use of paired measurements should be applied unless both autocorrelation and clinical circumstances indicate that the measurements are independent Check the presence of proportional bias, spread, or both, visually in the Bland–Altman plot and with regression analysis. If present, consider regression analysis to display the bias or LOA as a
Interpretation	Reference precision	function of the underlying CO, or data transformation Determine the repeatability of the experimental and reference technique for correct interpretation of the LOA and mean error
	Response time	Consider changes in CO and differences in response time if one or more continuous techniques disagree; if necessary and appropriate, measurements can be postponed
Data analysis	Trending ability	If applicable, consider the clinical concordance method as an alternative or addition to 4Q and polar analysis

by a visual check of the Bland–Altman plot. The spread of the differences may also be non-uniform over the range of CO measurements. This proportional spread can be identified visually in the Bland–Altman plot as a change in the scatter of the differences. In addition, the absolute values of the residuals as obtained with linear regression can be plotted against the mean CO.¹⁴ If the bias or LOA are non-uniform, transformation of the data or regression analysis can be applied to prevent under- and overestimation in specific measurement ranges;^{14 17} however, this limits the interpretation of the study results.

If the bias or LOA are uniform over the range of measurements, the difference between two systems is relatively larger in the lower range in comparison with the higher range. A uniform bias of 0.6 litre min⁻¹ represents a 20% mean deviation if CO is 3.0 litre min⁻¹, but a 10% deviation if CO is 6.0 litre min⁻¹. In contrast, if the bias or LOA are non-uniform, this percentage deviation may be constant. A non-uniform bias of 0.3 litre min⁻¹ at 3.0 litre min⁻¹ and of 0.6 litre min⁻¹ at 6.0 litre min⁻¹ represents a constant 10% deviation. Measurement error may therefore be constant in an absolute (e.g. 0.3 litre min⁻¹) or relative (e.g. 10%) sense.

Recommendation

Check the presence of proportional bias or spread visually in the Bland–Altman plot and with regression analysis. If present, consider regression analysis to display the bias or LOA as a function of the underlying CO, or data transformation.

Paired measurements

Many studies use multiple measurements in the same subject. Bland-Altman analysis without correction for paired measurements may underestimate the SD of the differences, leading to falsely narrow LOA and confidence intervals (CIs).^{5 6 12 14} As illustrated by Hamilton and Lewis,⁵ this effect increases with a small number of subjects, large number of measurements per subject, and little within-subject variance in comparison to betweensubject variance. This emphasizes the need for correction for paired measurements in studies investigating continuous CO monitoring devices in the absence of major haemodynamic changes. Consecutive measurements will tend to correlate, reducing the within-subject variance. In contrast, major haemodynamic changes may increase the within-subject variance to an extent that measurements become independent.¹⁸ We therefore suggest determining the autocorrelation of repeated measurements first. If this autocorrelation is not negligible, a correction for the use of paired measurements should be applied. Two methods are available for this purpose.⁶ ¹⁴ Bland and Altman¹⁴ provide a method to determine the LOA from the within-subject variances of the experimental and reference methods and the variance of the differences between the within-subject means. Alternatively, Myles and Cui⁶ use the average of repeated measurements and use a random effects model to correct for the reduction in variation that occurs by using this average. In addition to these statistical approaches, it is advisable to separate consecutive measurements in time, especially in the absence of major haemodynamic fluctuations. In this way, substantial correlation between consecutive measurements can be prevented.

Recommendation

A correction for the use of paired measurements should be applied unless both autocorrelation and clinical circumstances indicate that the measurements are independent. In the timing of consecutive measurements, the measurement protocol should consider the presence or absence of haemodynamic fluctuations.

Confidence intervals

Investigators should not forget to calculate 95% CIs for the bias, LOA and mean error, because they represent an estimation of their 'true' counterparts in a target population.⁷¹² At first sight, bias and LOA in a study may seem clinically acceptable. If, however, the corresponding CIs are wide, considerable differences between two systems can still be present in the target population. To illustrate this, we reconstructed the CIs of the bias, upper and lower LOA, and mean error in a number of studies (Supplementary Appendix A). Considering the CIs in the data analysis would probably lead to different conclusions in some studies. The CI of the bias should not be confused with the LOA.¹⁹ The CI of the bias indicates the limits for the bias in the target population, whereas the LOA refer to the spread of the differences in a specific study. The CI of the bias is calculated as $bias \pm t_{\alpha,n-1} + SD/\sqrt{n}$, and decreases with increasing sample size. Being a measure of spread, the LOA do not decrease by increasing the sample size.

Recommendation

The bias, LOA, and mean error should always be accompanied by their 95% CIs.

Agreement

Bland–Altman analysis does not provide definitive answers in terms of P-values. The acceptable level of agreement between a new and a reference CO technique is a matter of clinical judgment. A bias of 0.5 litre min⁻¹ and LOA of ±1.0 litre min⁻¹ may be acceptable for patients undergoing surgery with major haemodynamic disturbances, but not for patients with heart failure undergoing cardiac surgery. Clinically acceptable boundaries for bias and LOA or mean error should therefore always be defined in advance, depending on the target patients in which the new device is aimed to be used.³ To a certain extent, the desirable level of agreement can be adjusted if the new device has clear advantages over the reference technique in terms of safety, handling in clinical practice, or costs.

Recommendation

Acceptable boundaries for the bias, LOA, and mean error should be defined in advance.

Sample size calculations

The use of predefined criteria for Bland-Altman variables facilitates the decision-making process of accepting or rejecting new CO monitors for clinical use. However, study results may have the tendency to end up close to the predefined criteria, as these criteria reflect the clinical context in which the study has been performed. If the 95% CIs are wide, there is a substantial risk that they include the predefined criteria, which hinders definite conclusions. It is therefore advisable to consider the appropriate sample size in advance. Sample size calculations for Bland-Altman analysis can be considered controversial, because the method is not a statistical test. Moreover, the variability of (repeated) measurements with the new technique is unknown. Despite this, we point to a number of methods to estimate the appropriate sample size. First, the use of a desired maximal width for the 95% CIs around the mean error enables sample size calculations. This method was applied in a previous study by our group,²⁰ and is described in Supplementary Appendix B. Similar to this approach, the width of the CIs around the upper and lower LOA can be determined in terms of the sD, as described by Bland.²¹ Third, sample sizes can be estimated based on historical data. We realize that these approaches can be debated, and researchers are free to consider their use; however, we advise reflection on this topic in the design phase of method comparison studies in order to reduce the risk of underpowering.

Recommendation

Sample size calculations may be considered to estimate the appropriate number of subjects.

Reference precision

The LOA and mean error are influenced by the precision of the reference technique.⁴ ²² This is reflected in the formula by Critchley and Critchley²² to derive the mean error from the precision of the experimental and reference techniques, or:

Mean error =
$$\sqrt{([experimental precision]^2)}$$

+ [reference precision]²)

The use of imprecise reference techniques will therefore lead to wide LOA and high mean error, independent of the precision of the new device.⁸¹² Intermittent thermodilution CO (TDCO) with a pulmonary artery catheter is frequently used as reference technique. In many studies, the precision of TDCO is assumed to be 20%, and experimental precision should not exceed this 20% to be interchangeable with TDCO. Consequently, the mean error should not exceed $\sqrt{(20^2+20^2)}=28.3\%$, which is often rounded up to 30%.²² The strict use of a 30% limit for the mean error will, however, lead to erroneous conclusions if reference precision is significantly smaller or larger than 20%. Precision of TDCO or alternative techniques, such as transpulmonary thermodilution (TPCO), may even be improved to 5%.¹ ²³⁻²⁷ Both TDCO and TPCO can therefore be considered valuable as a reference technique, if properly performed. Moreover, this emphasizes the need for evaluation of reference precision in addition to experimental precision. The SD of repeated measurements or 'repeatability' can be used for this purpose.^{3 8} Repeatability is defined as 2×SD of repeated measurements (SD_{rep}) divided by CO.⁴ The squared values of experimental and reference repeatability can be added up as:

Combined repeatability = $\sqrt{([experimental repeatability]^2 + [reference repeatability]^2)}$

This 'combined repeatability' represents the maximal variation in repeated experimental and reference measurements that could explain the mean error. The mean error should therefore not exceed this value for the techniques to be interchangeable.³⁸

Recommendation

The TDCO and TPCO may be precise reference techniques, if properly performed. Both experimental and reference repeatability should be determined for proper interpretation of the LOA and mean error.

Changes in cardiac output and response time

Changes in CO introduce variability in repeated measurements, irrespective of precision (Fig. 2A). This does not affect the difference between experimental and reference CO if they are observed

at exactly the same moment in time (Fig. 2B). In the case of differences in response time between experimental and reference CO, however, a difference between the techniques will appear. This has important consequences for studies evaluating continuous devices during haemodynamic changes. These devices need time to process changes in the underlying CO, in contrast to intermittent reference techniques without measurement delay. Discrepancy will therefore emerge during haemodynamic changes, which fade out in time.²⁰ The timing and recording of measurements is therefore important, and postponing measurements during acute haemodynamic changes should be considered.²⁰ In acute settings, however, observations are directly followed by therapeutic decisions. To be valid in this situation, monitoring systems should display short response times.

Recommendation

The response time of (continuous) monitoring systems should be taken into account, and the method of collecting and recording CO measurements should be defined clearly. If necessary and appropriate, measurements can be postponed.

Trending ability

An increasing number of studies focus on the ability to track changes in CO, in addition to determining its absolute value.^{9 10} Evaluation of the trend in CO might be helpful to evaluate the effects of interventions and is intuitive, because CO is continuously changing as a result of a variety of influences, such as respiration, the autonomic nervous system, and changes in metabolic demand.^{24 25 28} The absolute value of CO is useful to consider in the diagnostic work-up of critical care patients. A proper evaluation of trending ability requires that changes in CO are induced in a controlled set-up. Moreover, the timing of and recording of measurements should be described clearly. Differences in response time between the experimental and reference method should be taken into account, and reference CO should be precise, as described earlier.

Bland–Altman analysis, polar plot methodology, and four-quadrant concordance

Although Bland–Altman analysis evaluates the accuracy and precision of absolute CO readings, conclusions about trending ability may be drawn intuitively. If absolute CO measurements are precise, trending ability should be adequate, irrespective of accuracy. Accuracy refers to the mean deviation between a new CO monitor and true CO. This deviation will be fixed in highly precise monitors and therefore irrelevant in tracking CO changes. In imprecise monitors, the deviation from the underlying CO is variable, which makes trending impossible. Theoretically, precision can be used to determine which changes in CO will be followed reliably. Precision of Δ CO is defined as $\sqrt{2}$ times the precision of a single CO measurement.⁴ As a result, a CO measurement device with a precision of 10% can reliably detect changes in CO of >14.1% ($\sqrt{2}$ *10). On the contrary, precision of CO measurement needs to be <7.1% (10/ $\sqrt{2}$) to detect a Δ CO of 10% reliably.

Two articles by Critchley and colleagues⁹ ¹⁰ review several methods to evaluate trending ability, including four-quadrant (4Q) concordance and polar plot methodology. The 4Q method plots the change in experimental CO (ΔCO_{exp}) against the change in reference CO (ΔCO_{ref}).⁹ The percentage of data points in which ΔCO_{exp} and ΔCO_{ref} change in the same direction is called 4Q concordance. This represents a rather crude estimate of trending ability and does not consider the magnitude of ΔCO_{exp} and ΔCO_{ref} .



Fig 2 The influence of changes in the underlying CO (A) and difference in response time between the experimental and reference technique (B). (A) Both variability in the underlying CO and imprecision may lead to significant spread in repeated measurements. The left box indicates variability in five consecutive CO measurements by a decrease in underlying CO, whereas imprecision in the measurement technique itself explains the variability in the measurements in the right box. (B) Changes in underlying CO in time (continuous line) do not affect the differences between the experimental (open crosses) and reference (filled crosses) CO measurements if the response time is nearly the same (left panel). If the experimental CO device responds slowly (dotted line, right panel), the differences increase, resulting in wide LOA. CO, cardiac output; LOA, limits of agreement.

In contrast, the polar plot approach enables quantitative assessment of trending ability, which is a major advantage.^{9 10} Nonetheless, a number of limitations need to be considered. First, interpretation of the polar variables is not straightforward. The translation of angular bias and radial LOA to clinical practice is not intuitive. Second, the criteria for good trending ability were validated, in a limited number of studies, against concordance and the opinion on trending ability by the authors. As a result, conclusions from polar plot analysis will have the tendency to agree with other statistical methods applied in the past, which limits the added value. Third, the criteria were determined with TDCO as the reference technique. In the case of another reference technique with different precision, the criteria should be adjusted.¹ Fourth, both polar plot and 4Q methods use exclusion zones to limit the influence of small changes in CO that may introduce random noise; however, this reduces statistical power and ignores potentially valuable information. The combination of small increases in ΔCO_{exp} (e.g. +1%) with small decreases in ΔCO_{ref} (e.g. -1%) or vice versa may be considered good trending, because these changes are both insignificant and unlikely to trigger therapeutic actions. In 4Q and polar analysis, these data pairs are excluded.

'Clinical concordance'

Alternatively, it is possible to pass a clinical judgment on each individual data pair. Each combination of ΔCO_{exp} and ΔCO_{ref} is

designated as 'good' or 'poor' trending and depicted in a 'clinical concordance' plot (Fig. 3A). The designations are based on criteria from a clinical perspective. Changes in CO_{ref} in a patient are divided into the following categories:

non-significant change ($\Delta CO_{ref} \pm 5\%$ or less);

moderate increase or decrease ($\Delta CO_{ref} \pm 5-15\%$); or large increase or decrease ($\Delta CO_{ref} \pm 15\%$ or more).

Each corresponding ΔCO_{exp} is assigned good trending if ΔCO_{exp} changes in the same direction and falls into the same category as ΔCO_{ref} . Depending on the clinical context, the number of categories and their criteria can be adjusted. 'Clinical concordance' can be defined simply as the percentage of 'good trending' assignments. The percentage 'poor trending' assignments directly informs the clinician about the risk for clinically relevant, erroneous trending information. Moreover, comparing the categories into which ΔCO_{exp} and ΔCO_{ref} fall provides insight into the extent to which ΔCO_{exp} and ΔCO_{ref} (dis)agree. In analogy with error grids used to validate new glucose measurement devices, this (dis)agreement can be further divided from the perspective of therapeutic consequences.²⁹ An error grid can be created to reflect the therapeutic consequences in specific zones in the concordance plot (Fig. 3B). The following zones can be distinguished.

(i) ΔCO_{exp} and ΔCO_{ref} change in the same direction and to the same extent, reflecting the following situations (in analogy with 'clinical concordance'): (a) CO_{exp} and CO_{ref} change



Fig 3 Clinical concordance and error grid plots. (A) Clinical concordance defines three categories (green squares), in which trending is 'good' from a clinical perspective. Clinical concordance represents the percentage of ΔCO_{exp} and ΔCO_{ref} data pairs falling into these categories. (B) The corresponding error grid uses multiple zones (rectangles in different shades of green) to define the level of agreement between ΔCO_{exp} and $\Delta \text{CO}_{\text{ref}}$ data pairs from the perspective of the rapeutic consequences. The zones are based on the clinical concordance categories, and can be created by extending the lines that surround the clinical concordance squares. Zone 1 corresponds to the clinical concordance categories in which CO_{exp} and CO_{ref} change in the same direction and to the same extent. This results in correct treatment decisions. In Zone 2, CO_{exp} and $\mathrm{CO}_{\mathrm{ref}}$ change in the same direction but not to the same extent, reflecting insufficient or exaggerated treatment. In Zone 3, CO_{exp} changes while CO_{ref} is constant or vice versa, reflecting unnecessary or withheld treatment. Zone 4 represents opposite changes in $\mathrm{CO}_{\mathrm{exp}}$ and $\mathrm{CO}_{\mathrm{ref}}$ resulting in opposite treatment. $\Delta \text{CO}_{\text{exp}}$, change in experimental cardiac output; CO_{ref}, change in reference cardiac output.

<5%; (b) CO_{exp} and CO_{ref} change between 5 and 15%; or (c) CO_{exp} and CO_{ref} change >15%. In this zone, correct treatment decisions are made with the new technique.

- (ii) ΔCO_{exp} and ΔCO_{ref} change in the same direction but not to the same extent, reflecting the following situations: (a) CO_{exp} changes between 5 and 15% while CO_{ref} changes >15%; or (b) CO_{exp} changes >15% while CO_{ref} changes between 5 and 15%. In this zone, treatment may be insufficient (a) or exaggerated (b).
- (iii) ΔCO_{exp} changes while ΔCO_{ref} is constant or vice versa, reflecting the following situations: (a) CO_{exp} changes >5% while CO_{ref} is constant; or (b) CO_{exp} is constant while CO_{ref} changes >5%. In this zone, unnecessary treatment may be initiated (a) or necessary treatment may be withheld (b).
- (iv) ΔCO_{exp} and ΔCO_{ref} change in opposite directions, reflecting the following situations: (a) CO_{exp} increases >5% while CO_{ref} decreases >5%; or (b) CO_{exp} decreases >5% while CO_{ref} increases >5%. In this zone, opposite treatment may be initiated.

The clinical concordance method provides a crude measure of trending agreement (clinical concordance) in combination with the therapeutic consequences of trending disagreements (error grid). A worked example is provided in Supplementary Appendix C. The suggested method uses all data pairs in the data analysis, which is an important advantage. Moreover, the extent to which ΔCO_{exp} and ΔCO_{ref} agree is addressed from a clinical perspective, which enhances the interpretation and use in clinical decision-making. The definitions for the clinical concordance categories and zones in the error grid are, however, rather subjective, and the use of different definitions might hinder comparison between studies in the future. Additional research is therefore needed to validate this new approach against current methods for trending analysis. Clinical concordance and error grids are meant as an extension to current methods, such as 4Q concordance and polar plot methodology, not as a substitute.

Recommendation

The clinical concordance method should be considered as an alternative or addition to 4Q and polar analysis in the evaluation of trending ability.

Discussion

The present review article describes the methodological challenges with the application of Bland–Altman and trending analysis in CO method comparison research. Moreover, the concept of clinical concordance and a corresponding error grid method is introduced to evaluate trending ability from a clinical perspective. Based on the items discussed, a stepwise approach to the design and data analysis of CO method comparison research can be created (Table 1). This approach may serve as a checklist for new researchers in the field. In addition, it may help clinicians to interpret the results from these studies in their decisions to incorporate new CO monitoring techniques in daily practice.

Although this review focuses on Bland–Altman and trending analysis, the data analysis of method comparison studies should not be restricted to these statistical methods. As in any type of research, the data analysis should include a close look at the raw data, considering outliers, haemodynamic circumstances, and patient characteristics. The scatterplot depicting experimental against reference CO should be evaluated, together with the range of CO measurements, effects in regions with high- or low-CO states, and effects in subgroups of patients. This is important because the performance of CO monitors may differ considerably depending on (patho)physiological conditions in the patient.^{1 2} Moreover, method comparison research represents only the initial part of the validation process of new CO monitors.³⁰ Besides technical efficacy, the ultimate goal of any newly developed monitor is to improve patient outcome and to be cost-effective. Method comparison studies should therefore anticipate application in clinical practice. This was an important reason to point to the use of predefined criteria defined within a desired, future clinical context. In addition, clinical concordance was introduced as a clinically intuitive method for trending ability, in which the level of (dis) agreement is translated to therapeutic consequences. Researchers should challenge themselves to embed the clinical context into their studies, both for a better understanding of the results and in order to facilitate the implementation of new technology in daily care.

Authors' contributions

All authors have been involved in drafting the manuscript, have given final approval of the version to be published, and agree to be accountable for all aspects of the work.

Supplementary material

Supplementary material is available at British Journal of Anaesthesia online.

Declaration of interest

W.F.B. has received honoraria for lectures and was consultant for Edwards Lifesciences. W.F.B. has no non-financial competing interests. The other authors declare that they have no competing interests.

Funding

Departmental resources.

References

- Peyton PJ, Chong SW. Minimally invasive measurement of cardiac output during surgery and critical care: a metaanalysis of accuracy and precision. Anesthesiology 2010; 113: 1220–35
- 2. De Waal EE, Wappler F, Buhre WF. Cardiac output monitoring. Curr Opin Anaesthesiol 2009; **22**: 71–7
- Mantha S, Roizen MF, Fleisher LA, Thisted R, Foss J. Comparing methods of clinical measurement: reporting standards for Bland and Altman analysis. Anesth Analg 2000; 90: 593–602
- Cecconi M, Rhodes A, Poloniecki J, Della Rocca G. Bench-tobedside review: the importance of the precision of the reference technique in method comparison studies – with specific reference to the measurement of cardiac output. Crit Care 2009; 13: 201–6
- Hamilton C, Lewis S. The importance of using the correct bounds on the Bland–Altman limits of agreement when multiple measurements are recorded per patient. J Clin Monit Comput 2010; 24: 173–5
- Myles PS, Cui J. Using the Bland–Altman method to measure agreement with repeated measures. Br J Anaesth 2007; 99: 309–11

- Hamilton C, Stamey J. Using Bland–Altman to assess agreement between two medical devices – don't forget the confidence intervals! J Clin Monit Comput 2007; 21: 331–3
- 8. Berthelsen PG, Nilsson LB. Researcher bias and generalization of results in bias and limits of agreement analyses: a commentary based on the review of 50 Acta Anaesthesiologica Scandinavica papers using the Altman–Bland approach. Acta Anaesthesiol Scand 2006; **50**: 1111–3
- 9. Critchley LA, Lee A, Ho AMH. A critical review of the ability of continuous cardiac output monitors to measure trends in cardiac output. Anesth Analg 2010; **111**: 1180–92
- Critchley LA, Yang XX, Lee A. Assessment of trending ability of cardiac output monitors by polar plot methodology. J Cardiothorac Vasc Anesth 2011; 25: 536–46
- Saugel B, Grothe O, Wagner JY. Tracking changes in cardiac output: statistical considerations on the 4-quadrant plot and the polar plot methodology. Anesth Analg 2015; 121: 514–24
- Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. Lancet 1986; 1: 307–10
- Bland JM, Altman DG. Comparing methods of measurement: why plotting difference against standard method is misleading. Lancet 1995; 346: 1085–7
- Bland JM, Altman DG. Measuring agreement in method comparison studies. Stat Methods Med Res 1999; 8: 135–60
- O' Brien E, Petrie J, Littler W, et al. The British Hypertension Society protocol for the evaluation of blood pressure measuring devices. J Hypertens 1993; 11: S43–62
- Razali N, Way YB. Power comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling tests. J Stat Model Analyt 2011; 2: 21–33
- Ludbrook J. Confidence in Altman–Bland plots: a critical review of the method of differences. Clin Exp Pharmacol Physiol 2010; 37: 143–9
- Jansen JR, Schreuder JJ, Settels JJ, et al. Single injection thermodilution. Anesthesiology 1996; 85: 481–90
- Nagele P. Misuse of standard error of the mean (SEM) when reporting variability of a sample. A critical evaluation of four anaesthesia journals. Br J Anaesth 2001; 90: 514–6
- 20. Montenij LJ, Buhre WF, De Jong SA, et al. Arterial pressure waveform analysis versus thermodilution cardiac output measurement during open abdominal aortic aneurysm repair: a prospective, observational study. *Eur J Anaesthesiol* 2015; **32**: 13–9
- Bland M. How can I decide the sample size for a study of agreement between two methods of measurement? Available from http://www-users.york.ac.uk/~mb55/meas/ sizemeth.htm (accessed 13 June 2015)
- 22. Critchley LA, Critchley JA. A meta-analysis of studies using bias and precision statistics to compare cardiac output measurement techniques. J Clin Monit Comput 1999; **15**: 85–91
- Jansen JR, Vesprille A. Improvement of cardiac output estimation by the thermodilution method during mechanical ventilation. Intensive Care Med 1986; 12: 71–9
- 24. Nishikawa T, Dohi S. Errors in the measurement of cardiac output by thermodilution. *Can J Anaesth* 1993; **40**: 142–53
- Jansen JR, Schreuder JJ, Settels JJ, Kloek JJ, Versprille A. An adequate strategy for the thermodilution technique in patients during mechanical ventilation. Intensive Care Med 1990; 16: 422–5
- Monnet X, Persichini R, Ktari M, et al. Precision of the transpulmonary thermodilution measurements. Crit Care 2011; 15: R204

- Gondos T, Marjanek Z, Kisvarga Z, et al. Precision of transpulmonary thermodilution: how many measurements are necessary? Eur J Anaesthesiol 2009; 26: 508–12
- Michard F, Teboul JL. Using heart-lung interactions to assess fluid responsiveness during mechanical ventilation. Crit Care 2000; 4: 282–9
- 29. Clarke WL, Cox D, Gonder-Frederick LA, Carter W, Pohl SL. Evaluating clinical accuracy of systems for self-monitoring of blood glucose. *Diabetes Care* 1987; **10**: 622–8
- 30. Pearl WS. A hierarchical outcome approach to test assessment. Ann Emerg Med 1999; **33**: 77–84

Handling editor: J. G. Hardman