

Effects of observation errors in linear regression and bin-average analyses

By H. L. TOLMAN*

National Centers for Environmental Prediction, USA

(Received 7 June 1996; revised 21 January 1997)

SUMMARY

Effects of observation errors in linear regression and bin-averaged (BA) validation techniques are investigated using the example of marine wind speeds. It is shown that a conventional linear regression systematically underestimates the slope of the regression line, and systematically overestimates the random model error. A BA analysis systematically underestimates extreme wind speeds, incorporates spurious nonlinearity, and overestimates random model errors. Correction techniques are suggested for studies in which the observation error can be estimated. Using synthetic data the potential of the correction techniques is illustrated, and it is shown that the above errors are generally not negligible for wind speed validation studies. Practical examples consider the random errors of anemometers and wind speed estimates from satellites. These examples highlight the importance of the error corrections, and illustrate the difficulty of estimating observation errors. Finally, it is argued that the well-known symmetric slope regression should not be used for the validation of forecast systems. Although the present study deals with marine wind speeds, its results are expected to be valid for a wide range of validation studies.

KEYWORDS: Observation error Regression Statistical methods

1. INTRODUCTION

In atmospheric and oceanic sciences linear regression analyses are used extensively to validate models and retrieval algorithms for remotely sensed data. It is well known that a conventional linear regression analysis is valid only if the observation error is negligible† (e.g. Draper and Smith (1981), section 2.14), but even in textbooks effects of observation errors are rarely discussed in detail. Whereas some validation studies attempt to account for observation errors by using more advanced linear regression techniques (see later), many validation studies simply ignore observation errors.

By definition, linear regression assumes quasi-linear model behaviour. Nonlinear model behaviour is often investigated using bin-averaged (BA) analyses, where model statistics are determined for several narrow ranges of observations. As will be shown, BA analyses are also sensitive to observation errors.

The present study addresses effects of observation errors in the above analysis techniques, and suggests error correction methods for studies where the observation error can be estimated independently. For convenience of discussion the validation of marine wind speeds with buoy observations is considered. The results are nevertheless expected to be valid for a wide range of validation studies in a wide range of disciplines. The analysis methods and effects of observation errors are discussed in section 2. Correction methods are described in section 3 and are tested with synthetic data in section 4. In section 5 two practical examples are presented. In the first, random anemometer errors are estimated. In the second, biases for satellite wind speeds are estimated. These applications illustrate the impact of the suggested error corrections, and difficulties in estimating observation errors. The advantages and limitations of the error-corrected analysis methods are discussed in section 6 and the conclusions are given in section 7.

* Corresponding address: Ocean Modeling Branch, Environmental Modeling Center, NOAA/NCEP, 5200 Auth Road, Room 209, Camp Springs, MD 20746, USA.

† With the exception of ‘controlled’ observations (Berkson 1950), which usually cannot be obtained in oceanography or meteorology (e.g. Ricker 1973).

2. ANALYSIS METHODS

Systematic errors in statistical analysis techniques are most easily investigated using continuous probability density functions (pdfs), and integral quantities of pdfs as discussed in subsection 2(a). Linear regression analyses and their errors are discussed in subsection 2(b), and BA analyses are discussed in subsection 2(c).

(a) Data description

The present study deals with errors in the estimation of errors. To minimize the inherent confusion a systematic notation is adopted. Upper-case symbols denote true (error-free) values, lower-case symbols denote estimates (including effects of errors), and a prime identifies an error quantity. The suffices ‘m’ and ‘o’ denote model and observation, respectively. For instance, defining σ as a standard deviation, Σ is the true standard deviation of the winds, Σ'_m is the true standard deviation of the model error and σ'_m is its estimate.

Consider a true wind speed U with a given pdf $P(U)$. Modelled and observed wind speeds generally contain systematic and random errors. For observations to be useful in a validation study their systematic errors should be negligible and/or removed from the data before the validation takes place ($U_o \equiv U$). Systematic model errors are to be retrieved, and thus are part of the true model wind speed $U_m \neq U$.

Model and observation errors are expected to be uncorrelated. Describing the observation error for a given wind speed U with the error pdf $p'_o(u_o|U)$, and using a similar description for the model error, the joint pdf $p(u_o, u_m)$ of observed and modelled wind speeds becomes

$$p(u_o, u_m) = \int_0^\infty P(U) p'_o(u_o|U) p'_m(u_m|U) dU . \tag{1}$$

This distribution is generally described with integral quantities such as the mean observed and modelled speeds \bar{u}_o and \bar{u}_m , the variances s_{oo} and s_{mm} and the covariance s_{om} ,

$$s_{oo} = \overline{u_o^2} - \bar{u}_o^2 , \tag{2}$$

$$s_{mm} = \overline{u_m^2} - \bar{u}_m^2 , \tag{3}$$

$$s_{om} = \overline{u_o u_m} - \bar{u}_o \bar{u}_m . \tag{4}$$

These parameters incorporate effects of random errors and, for random errors with symmetric distributions, they can be expressed in terms of their expected error-free counterparts as

$$\bar{u}_o = \langle U_o \rangle = \langle U \rangle , \tag{5}$$

$$\bar{u}_m = \langle U_m \rangle , \tag{6}$$

$$s_{oo} = S_{oo} + \langle S'_{oo} \rangle = S + \langle S'_{oo} \rangle , \tag{7}$$

$$s_{mm} = S_{mm} + \langle S'_{mm} \rangle , \tag{8}$$

$$s_{om} = S_{om} , \tag{9}$$

where $\langle \dots \rangle$ represents the expected value

$$\langle X \rangle = \int_0^\infty X(U) P(U) dU . \tag{10}$$

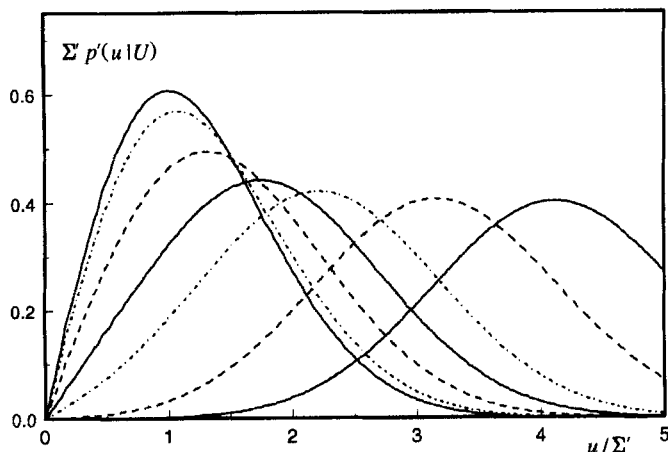


Figure 1. Normalized error probability density functions $\Sigma' p'(u|U)$ of Eq. (15) as a function of the normalized wind speed u/Σ' for true wind speeds $U/\Sigma' = 0, 0.5, 1, 1.5, 2, 3$ and 4 from left to right. The left-most curve represents the Rayleigh distribution (11) with $u_r = \Sigma'$. See text for further explanation.

The mean wind speeds \bar{u}_o and \bar{u}_m and the covariance s_{om} thus are not influenced by random errors, but the variances s_{oo} and s_{mm} are increased by the corresponding mean error variances $\langle S'_{oo} \rangle$ and $\langle S'_{mm} \rangle$. Note that the equalities (5) to (9) are valid here because the averages on the left side are calculated from the continuous pdfs. For practical studies these averages are calculated from the data and hence will include sampling errors, making the above equalities approximations.

To quantify errors in the following sections, pdfs have to be assumed. Somewhat arbitrarily, skewed wind speed distributions will be described here with a Rayleigh distribution (see Fig. 1)

$$P(U) = \frac{U}{u_r^2} \exp\left(\frac{-U^2}{2u_r^2}\right), \quad (11)$$

where u_r is the single parameter defining the distribution. The corresponding mean wind speed $\langle U \rangle$ and standard deviation Σ are

$$\langle U \rangle = \sqrt{\frac{\pi}{2}} u_r, \quad \Sigma = \sqrt{\frac{4-\pi}{2}} u_r. \quad (12)$$

All further calculations have been performed with $u_r = 7 \text{ m s}^{-1}$, resulting in $\langle U \rangle = 8.8 \text{ m s}^{-1}$ and $\Sigma = 4.6 \text{ m s}^{-1}$.

Random errors are commonly described with a normal distribution

$$p'(u|U) = \frac{1}{\Sigma' \sqrt{2\pi}} \exp\left(\frac{-(u-U)^2}{2\Sigma'^2}\right). \quad (13)$$

This distribution, however, includes negative wind speeds, which is obviously not realistic. Negative wind speeds can be avoided by considering the vector equivalent of (13)

$$p'(u|U) = \frac{1}{\Sigma'^2 2\pi} \exp\left(\frac{-\|\mathbf{u} - \mathbf{U}\|^2}{2\Sigma'^2}\right). \quad (14)$$

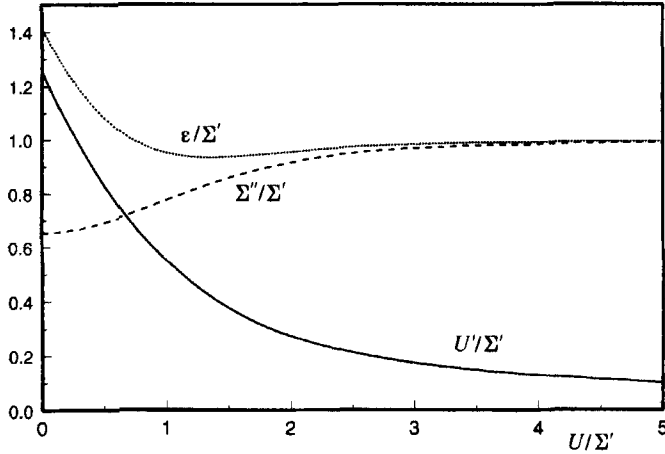


Figure 2. Normalized pseudo bias U'/Σ' , standard deviation Σ''/Σ' and r.m.s. error ϵ/Σ' as functions of the wind speed U/Σ' for the cut-off normal distribution (15). See text for explanation of the symbols.

Integration for constant $\|U\|$ gives the corresponding scalar pdf

$$p'(u|U) = \frac{u}{\Sigma'^2} I_0\left(\frac{uU}{\Sigma'^2}\right) \exp\left(-\frac{(u^2 + U^2)}{2\Sigma'^2}\right), \tag{15}$$

where $I_0(\dots)$ is a modified Bessel function (see, for example, Abramowitz and Stegun (1973), section 9.6.16). This pdf is presented in Fig. 1. For large U/Σ' , (15) corresponds to the normal distribution (13). For $U = 0$, (15) corresponds to the Rayleigh distribution (11).

Distribution (15) includes a pseudo bias U' (solid line in Fig. 2). This pseudo bias occurs for any assumed error distribution (e.g. Hinton and Wylie 1985), and makes it virtually impossible to separate mean and random errors for small wind speeds U/Σ' . Furthermore, the actual standard deviation Σ'' of (15) is considerably smaller than Σ' for wind speeds $U/\Sigma' < 3$ (dashed line in Fig. 2). The root-mean-square error $\epsilon = \sqrt{U'^2 + \Sigma''^2}$ differs by less than 10% from Σ' for $U/\Sigma' > 0.5$. Because the skewed errors are generally relevant for a relatively small part of the data only, and because the corresponding pseudo biases can be removed, symmetric errors are assumed in the following.

(b) *Linear regression*

In a linear regression analysis the functional relation $U_m = \Phi(U_o)$ is approximated by a straight line

$$u_m = a + bu_o, \tag{16}$$

where b is the regression coefficient and a is the intercept. Regression techniques differ in the way in which this relation is fitted to the data. The ‘conventional’ regression is the regression of u_m on u_o , where the optimization considers differences between the regression line and the modelled wind speeds only. Using the suffix ‘o’ for this regression line, the standard result becomes

$$a_o = \bar{u}_m - b_o \bar{u}_o, \quad b_o = \frac{s_{om}}{s_{oo}}. \tag{17}$$

Its error-free counterpart is

$$A_o = \langle U_m \rangle - B_o \langle U_o \rangle, \quad B_o = \frac{S_{om}}{S_{oo}}. \quad (18)$$

The latter regression line represents a better estimate of the functional relation $U_m = \Phi(U_o)$ as it is not influenced by random errors. The two regression lines have the point (\bar{u}_o, \bar{u}_m) in common, but have different slopes (regression coefficients). The effect of random errors can be estimated by normalizing b_o with B_o using (7) and (9)

$$\frac{b_o}{B_o} = \left(1 + \frac{\langle S'_{oo} \rangle}{S_{oo}} \right)^{-1}. \quad (19)$$

The regression coefficient b_o thus systematically underestimates its error-free value B_o if the expected random observation error $\langle S'_{oo} \rangle$ is not negligible (see also, for example, Draper and Smith (1981), section 2.14).

An alternative regression of u_o on u_m (denoted here as the inverse regression) is determined by considering the differences between the regression line and the observations only. Using the suffix 'm', the standard results and the normalized regression coefficient become, using (8) and (9),

$$a_m = \bar{u}_m - b_m \bar{u}_o, \quad b_m = \frac{S_{mm}}{S_{om}}, \quad (20)$$

$$\frac{b_m}{B_m} = 1 + \frac{\langle S'_{mm} \rangle}{S_{mm}}. \quad (21)$$

The regression coefficient b_m thus is systematically overestimated due to the occurrence of a random model error $\langle S'_{mm} \rangle$. Note that it is assumed that $B_o \approx B_m$, which is valid for quasi-linear relations.

More advanced linear regression techniques generally require that the ratio of random model and observation errors is known (e.g. Lindley 1947; Jolliffe 1990). The regressions (17) and (20) represent the limiting cases for the slope in the case of dominant random model or observation errors, respectively. As a representative of more advanced techniques, the geometric mean regression b_{gm} is considered

$$b_{gm} = \sqrt{\frac{S_{mm}}{S_{oo}}}, \quad (22)$$

which implies that the normalized random model and observation errors are identical, and which is also known as the 'symmetric slope'. Its normalized regression coefficient follows from Eq. (18) as

$$\frac{b_{gm}}{B_o} \approx \sqrt{\frac{1 + \langle S'_{mm} \rangle / S_{mm}}{1 + \langle S'_{oo} \rangle / S_{oo}}}. \quad (23)$$

The geometric mean regression thus corresponds to the true functional regression B_o for $\langle S'_{oo} \rangle / \langle S'_{mm} \rangle = S_{mm} / S_{oo}$. If the random observation error is larger (smaller) b_{gm} underestimates (overestimates) B_o .

Apart from the functional relation Φ , random model errors $\langle S'_{mm} \rangle \approx \bar{s}'_{mm}$ are assessed in validation studies. In the conventional regression (17) the model error is estimated as

$$\overline{s'_{mm}} = s_{mm} - \frac{s_{om}^2}{s_{oo}} = s_{mm} - b_o s_{om} . \tag{24}$$

Generally, the model error is estimated from any regression coefficient b as

$$\overline{s'_{mm}} = s_{mm} - b s_{om} . \tag{25}$$

Thus the model error is overestimated (underestimated) if b underestimates (overestimates) B_o . For the limiting cases the entire observation error is erroneously attributed to the model (b_o), or the model error is by definition negligible (b_m). In the geometric mean regression the error is equally distributed over the model and the observations.

(c) *Bin-average (BA) analysis*

Determining statistics per class of observed speeds u_o corresponds to calculating n th-order moments $m_n(u_o)$ from a continuous pdf

$$m_n(u_o) = \int_0^\infty u_m^n p(u_o, u_m) du_m , \tag{26}$$

for given data intervals or ‘bins’ δu_o . For the moment, effects of a finite bin size will be ignored. The moments m_0 through m_2 yield the marginal pdf of the observed speed $p_o(u_o)$, the mean modelled speed $\bar{u}_m(u_o)$, and the error variance $s'_{mm}|u_o$

$$m_0 = p_o(u_o) , \tag{27}$$

$$m_1 = p_o(u_o) \bar{u}_m(u_o) , \tag{28}$$

$$m_2 = p_o(u_o) (s'_{mm}|u_o + \bar{u}_m^2(u_o)) . \tag{29}$$

For convenience, the explicit dependence on u_o will henceforth be dropped from the notation. Substitution of (1) in (26) and rearranging the order of integration yields

$$m_n = \int_0^\infty p'_o(u_o|U) M_n(U) dU , \tag{30}$$

where

$$M_n(U) = \int_0^\infty u_m^n P(U) p'_m(u_m|U) du_m , \tag{31}$$

represents the true moments of the model wind speed. The moments M_0 through M_2 yield the distribution P , mean model speed U_m and model error S'_{mm} as in Eqs. (27) to (29). Equation (30) shows that the observed moments are a convolution of the observation error $p'_o(u_o|U)$ and the true moments.

To assess the magnitude of errors introduced by a BA analysis, Eq. (30) has been evaluated for a ‘perfect’ model ($U_m = U_o = U$). Following common practice, errors are described as a standard deviation $\sigma'_o = \sqrt{s'_{oo}}$ etc., rather than as a variance. Arbitrarily, the model errors are set to be $\Sigma'_m = 1.5 \text{ m s}^{-1}$ ($\approx 0.17\langle U \rangle$), and the observation error is assumed to be a fraction of the wind speed, $\Sigma'_o = \alpha U$. The resulting normalized errors in the wind speed distribution $(p_o - P)/P$, the functional behaviour $(\bar{u}_m - U_m)/\langle U \rangle$, and the random model error $(\sigma'_m - \Sigma'_m)/\Sigma'_m$ are presented as a function of the normalized wind speed $U_o/\langle U \rangle$ in Fig. 3.

The BA analysis reproduces the expected results if no observation errors are present ($\alpha = 0$). Note that the latter results in panels (b) and (c) incorporate effects of the skewed

error distribution for small wind speeds (cf. Fig. 2). Systematic errors of the BA analysis are identified as the difference between results for $\alpha \neq 0$ and $\alpha = 0$. Figure 3(a) shows that the BA analysis predicts the distribution P accurately, except for extreme wind speeds ($U/\langle U \rangle > 2.5$), where P is severely overestimated. Figure 3(b) shows similar results for the estimated model wind speed \bar{u}_m , which is systematically underestimated for high wind speeds ($U_o/\langle U \rangle > 2$). Figure 3(c) shows the BA analysis systematically overestimates the model error Σ'_m . This could be expected as the observation error is erroneously attributed to the model. Unlike in panels (a) and (b), this analysis error also depends on the random model error, as the second moment (29) includes the model error. The results of Fig. 3 are obviously sensitive to the choice of the observation error. In particular, if the observation error is finite for $U \rightarrow 0$, errors in p_o , \bar{u}_m and σ'_m will become larger for small wind speeds. However, they generally remain much smaller than the present errors for large $U/\langle U \rangle$ (figures not presented here).

As mentioned above, moments (26) are estimated from practical data for a finite bin width δu_o

$$m_n(u_o) \approx \frac{1}{\delta u_o} \int_{u_o - \delta u_o/2}^{u_o + \delta u_o/2} m_n(x) dx. \quad (32)$$

Errors introduced by this averaging are illustrated in Fig. 4 for several bin sizes $\delta u_o = \beta \langle U \rangle$. Errors are identified as the deviation of results for finite bin sizes ($\beta > 0$, symbols) from the 'exact' results ($\beta = 0$, solid lines).

If the bin size becomes too large to describe the data distribution adequately, the averaging over the bin is bound to result in noticeable errors. This explains the occurrence of errors in all panels of Fig. 4 for the largest values of β . Note that many meteorological studies use such large bin sizes by considering a very small number of bins (for instance, low, medium and high wind speeds only). If the bin size is sufficiently small to resolve the distribution of the data, p_o , \bar{u}_m and s'_{mm} in Eqs. (27) to (29) can be linearized in the integral (32). It is then easily shown that p_o and \bar{u}_m are not influenced by the finite bin size, but that approximately $\frac{1}{12} \delta u_o^2$ is added to the estimate of the model variance s'_{mm} . This will only be relevant if the bin size δu_o becomes larger than the local model error Σ'_m .

3. ERROR CORRECTION

The previous section shows that conventional linear regression and BA analyses incorporate systematic errors if the observation error is not negligible. If the observation error $\sigma'_o \approx \Sigma'_o$ is known or can be estimated, it is possible to estimate and correct such errors.

Using Eqs. (18), (9) and (7) the exact regression coefficient B_o can be estimated as

$$B_o = \frac{s_{om}}{s_{oo} - \langle S'_{oo} \rangle} \approx \frac{s_{om}}{s_{oo} - s'_{oo}} = b_{o,c}, \quad (33)$$

which can be interpreted as an error-corrected version of the conventional regression (17). The corresponding corrected estimate of the model error \bar{s}'_{mm} is obtained from Eq. (25). Approximating the true wind speed U and distribution P with its estimates u_o and p_o , the mean observation error \bar{s}'_{oo} in (33) can be estimated as

$$\bar{s}'_{oo} = \frac{1}{n} \sum_{i=1}^n \{ \sigma'_o(u_{o,i}) \}^2, \quad (34)$$

where n is the number of observations and $u_{o,i}$ represents individual observations.

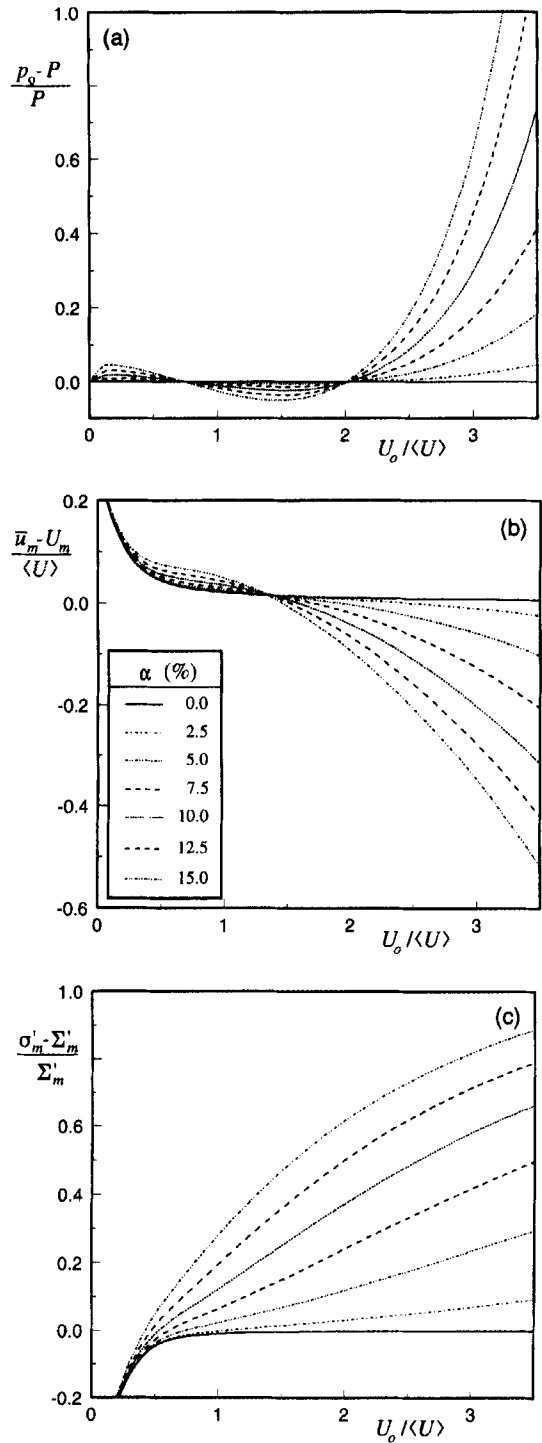


Figure 3. Normalized results of a bin-averaged analysis as a function of the normalized wind speed $U_o/\langle U \rangle$ for several observation errors $\Sigma'_o = \alpha U$. (a) Wind speed distributions ρ_o , (b) mean model wind speeds \bar{u}_m , and (c) random model error σ'_m . A 'perfect' model ($U = U_o = U_m$) with a random error $\Sigma'_m = 1.5 \text{ m s}^{-1} = 0.17(U)$. See text for further explanation.

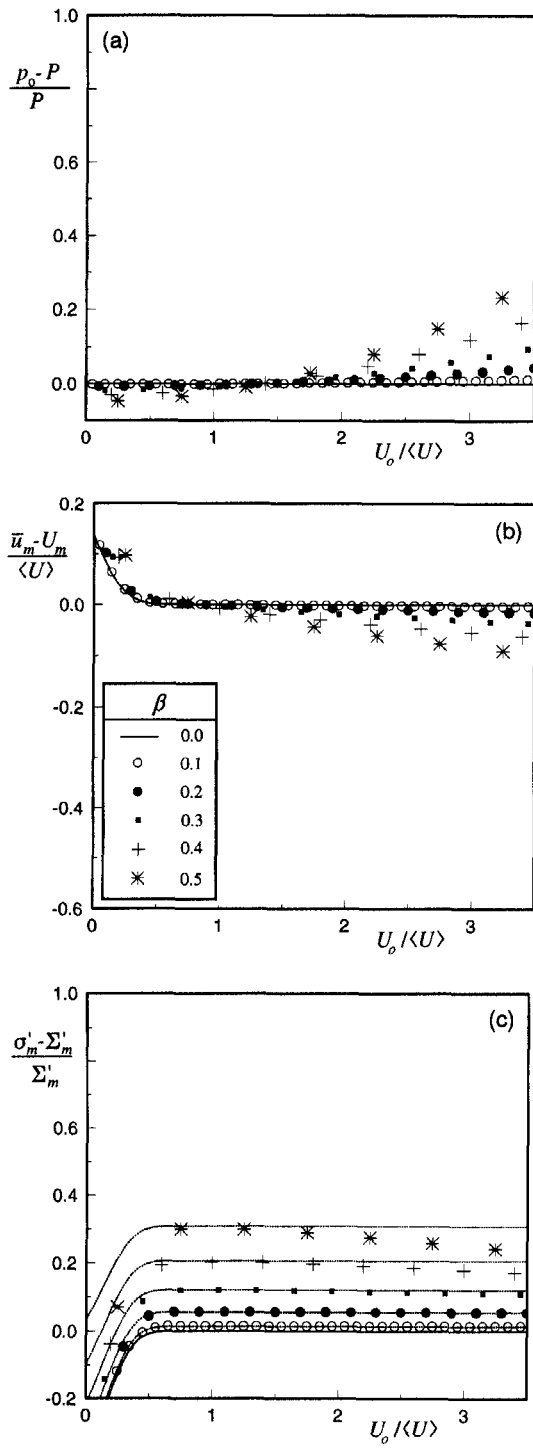


Figure 4. Like Fig. 3 but for a bin size $\delta u_o = \beta \langle U \rangle$. No observation error ($\alpha = 0$). Dotted line in (c): $\frac{1}{12} \delta u_o^2$ added to true model error variance.

The results of a BA analysis (denoted as $\bar{u}_{m,ba} = \phi_{ba}$ and $\sigma'_{m,ba}$, respectively) incorporate analysis errors due to the observation error (Fig. 3) and the finite bin size (Fig. 4). For the analysis to be useful, the bin size δu_o has to be sufficiently small to describe the data distribution accurately. The bin size then does not effect ϕ_{ba} , but systematically increases $\sigma'_{m,ba}$. A corrected estimate $\sigma'_{m,bc}$ can be obtained as

$$\sigma'_{m,bc} = \sqrt{(\sigma'_{m,ba})^2 - \frac{1}{12} \delta u_o^2}. \tag{35}$$

Estimating effects of observation errors requires evaluation of the moments (26) or (30) of the joint pdf. Again approximating U and P with u_o and p_o , estimates \tilde{m}_0, \tilde{m}_1 and \tilde{m}_2 for the moments m_0, m_1 and m_2 can be calculated as

$$\tilde{m}_0 = \frac{1}{n} \sum_{i=1}^n p'_o(u_o|u_{o,i}), \tag{36}$$

$$\tilde{m}_1 = \frac{1}{n} \sum_{i=1}^n \tilde{\phi} p'_o(u_o|u_{o,i}), \tag{37}$$

$$\tilde{m}_2 = \frac{1}{n} \sum_{i=1}^n (\tilde{\phi}^2 + (\tilde{\sigma}'_m)^2) p'_o(u_o|u_{o,i}), \tag{38}$$

where $\tilde{\phi}$ and $\tilde{\sigma}'_m$ are estimates of the mean model behaviour and the random model error, respectively (evaluated at $u_{o,i}$). The corresponding errors $\Delta\tilde{\phi}$ and $\Delta\tilde{\sigma}'_m$ then can be estimated as

$$\Delta\tilde{\phi} = \frac{\tilde{m}_1}{\tilde{m}_0} - \tilde{\phi}. \tag{39}$$

$$\Delta\tilde{\sigma}'_m = \left\{ \frac{\tilde{m}_2}{\tilde{m}_0} - \left(\frac{\tilde{m}_1}{\tilde{m}_0} \right)^2 \right\}^{1/2} - \tilde{\sigma}'_m. \tag{40}$$

Obvious first guesses for $\tilde{\phi}$ and $\tilde{\sigma}'_m$ would be to assume that the mean model behaviour is perfect, and that the random model error can be estimated using Eq. (35)

$$\tilde{\phi}(u_{o,i}) = u_{o,i}, \tag{41}$$

$$\tilde{\sigma}'_m(u_{o,i}) = \sigma'_{m,bc}(u_{o,i}), \tag{42}$$

where $\sigma'_{m,bc}$ needs to be described with a polynomial (or other) fit to extend it over the entire data range. Because both first guesses (41) and (42) potentially include significant errors, it is prudent to repeat the estimation of the analysis errors with estimates of the model behavior that include the error estimates (41) and (42). This estimate of the analysis error will be denoted as the second guess. Although additional iterations are easily made, they are expected generally to have a limited impact on the final results.

An additional error occurs in this analysis because the observed distribution $p_o(u_o)$ will be broader than the true distribution $P(U)$ because of its convolution with the observation error. This effectively stretches the u_o axis, and can be corrected by statistically adjusting the observed wind speeds. The necessary shift of the observed wind speed is the difference between the true and the expected observed wind speeds, which approximately

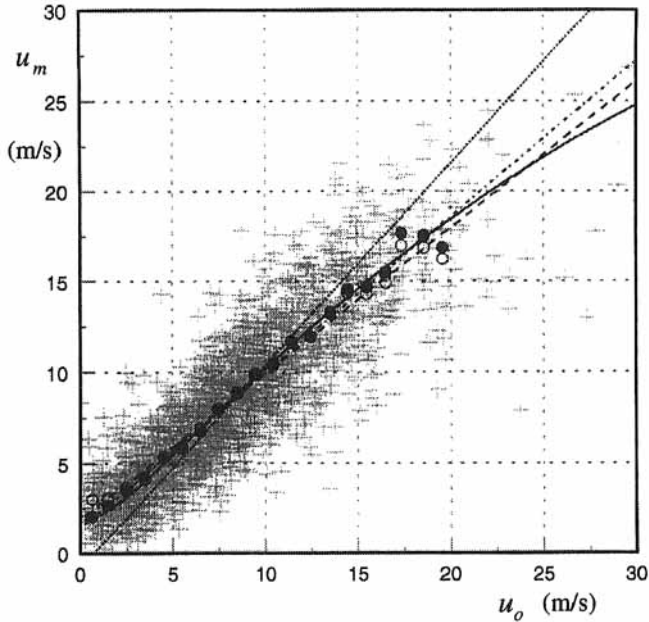


Figure 5. Monte Carlo realization of measured and modelled wind speeds u_o and u_m , respectively. + (simulated data); solid line (true model behaviour $U_m = \Phi(U)$); dashed line (conventional regression (Eq. (17))); dotted line (inverse regression (Eq. (20))); chain line (Error-corrected regression (Eqs. (33) and (34))); \circ (bin-averaged analysis for a bin width of 1 m s^{-1} and requiring a minimum of 10 observations per bin); and \bullet (corresponding error-corrected bin-averaged analysis).

equals the above first-guess correction $\Delta\tilde{\phi}$. This error correction is thus most elegantly included in the second-guess correction, and automatically accounts for pseudo biases in the observations for low wind speeds as shown in Fig. 2.

Henceforth, the error corrected BA method will be denoted as the ECBA method. Second-guess corrections and third-order polynomial descriptions of ϕ_{ba} and $\sigma'_{m,bc}$ are used unless specified differently.

4. MONTE CARLO EXPERIMENTS

To illustrate the analysis errors and the potential of the correction techniques, Monte Carlo experiments have been performed. The observations are assumed to be free of bias ($U_o = U$), and the random observation error $\Sigma'_o = \max(1, 0.1U)$ is assumed to be known. This error corresponds to the required system accuracy of wind observations from buoys (Gilhausen 1987) (actual observation errors are discussed in the following section). To simulate nonlinear model behaviour the model wind speed is defined as $U_m = 1.1U - 0.01U^2$. Arbitrarily, the random model error is defined as $\Sigma'_m = \max(1.5, 0.2U)$, with a smooth transition between the two branches. The sample size is set to 2000, which corresponds to three months of hourly observations at a single location. A complete Monte Carlo experiment consists of many realizations of this model. However, as it is used here only as an illustration, the discussion will focus on a single realization of this model (Fig. 5). Deviations from other realizations will be mentioned where necessary. To avoid unnecessary complications, the pseudo bias for low model wind speeds (solid line in Figs. 2 and 3(b)) is included in U_m (as is generally the case in practical studies).

TABLE 1. REGRESSION COEFFICIENTS (b) CALCULATED FROM THE EQUATIONS INDICATED, AND MODEL ERRORS ($\overline{\sigma'_m}$, $\langle \Sigma'_m \rangle$ EQ. (25)) FOR THE MONTE CARLO EXPERIMENT OF FIG. 5

	b	Equation number	$\overline{\sigma'_m}$, $\langle \Sigma'_m \rangle$ (m s^{-1})
Conventional	0.806	(17)	2.32
Corrected conventional	0.860	(33)	2.11
Exact	0.891	(18)	2.12
Geometric mean	0.950	(22)	1.70
Inverse	1.119	(20)	0

Results of several regression analyses are shown in Fig. 5 and Table 1. The conventional (b_o) and inverse regression (b_m) incorporate the largest errors, and establish upper and lower boundaries of analysis errors. The geometric mean regression (b_{gm} , not in figure) also includes significant errors as the model error is much larger than the observation error. The best results are obtained with the error-corrected conventional regression ($b_{o,c}$), which underestimates the regression coefficient by 3% and reproduces Σ'_m with a negligible error. The error in $b_{o,c}$ is anomalously large compared with other Monte Carlo realizations, which generally show analysis errors of less than 1%. The error-corrected conventional regression thus is virtually free of errors for most Monte Carlo realizations. It does not, however, describe the true nonlinear model behaviour accurately, in particular for extreme wind speeds (compare chain and solid lines in Fig. 5.)

A BA analysis requires the choice of a bin width and a required minimum number of data per bin. In the present experiment a small bin width (1 m s^{-1}) and a small minimum number of observations per bin (10 for $\bar{u}_m = \phi$ and 20 for σ'_o) are chosen to illustrate both effects of sampling errors and the potential of the (corrected) BA method. For wind speeds below 15 m s^{-1} the results show little effects of sampling variability due to the relatively large number of data per bin (30 to 200). For wind speeds over 15 m s^{-1} the number of data per bin drops quickly, and the corresponding results display an increasing sampling error.

Mean model errors are isolated as residuals $u_m - U$ in Fig. 6(a). The BA analysis overestimates the true model wind speed for low winds, and underestimates it for high wind speeds. For high wind speeds this is somewhat obscured by the sampling errors, but it becomes evident if a third-order polynomial is fitted to the results. The ECBA results display much smaller analysis errors, and the deviation from the exact solution appears to be dominated by sampling errors. Estimated analysis errors are presented in Fig. 6(b). Both the first and second guess show some oscillations but give a fair representation of the expected error (calculated directly from the known joint distribution). The oscillations for $u_o > 20 \text{ m s}^{-1}$ appear to be related to the sparsity of the data in this wind regime, and were found to be significantly larger for other Monte Carlo realizations. As the analysis does not render results in this regime, these oscillations are not relevant for the error correction. Differences between the first and second guess are small, because effects of the improved estimate of the mean model behaviour and the correction of the observations mostly cancel. This is not generally the case.

Estimated random model errors σ'_m are presented in Fig. 7(a). The results of the BA analysis overestimate the model error Σ'_m systematically. The ECBA analysis again shows a significant improvement. Results for the highest wind speeds tend to deviate more from the true model behaviour than ϕ , because the random error is more sensitive to sampling than the mean behaviour. In other Monte Carlo realizations, model errors Σ'_m for

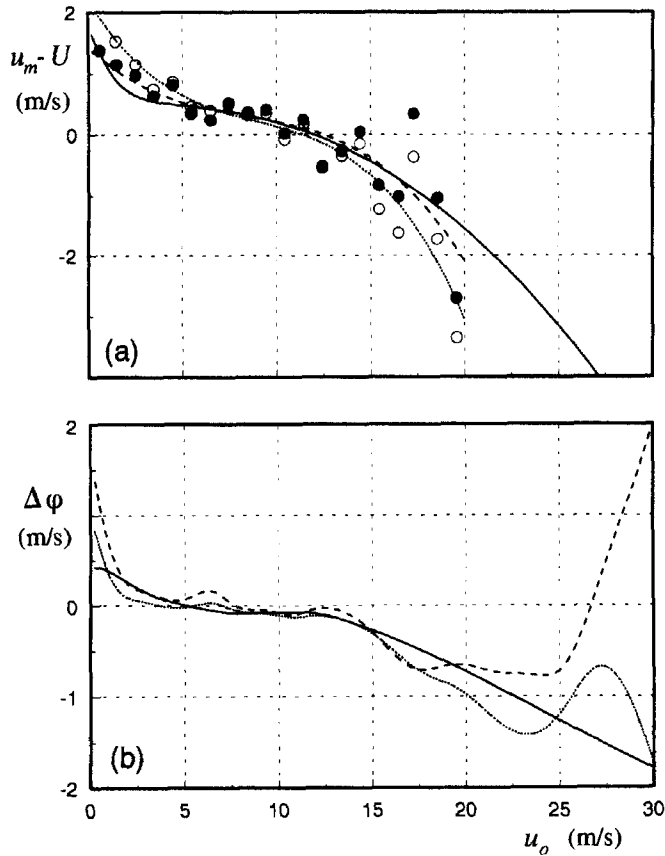


Figure 6. (a) Residual wind speeds $u_m - U$ as a function of the observed wind speed u_o corresponding to Fig. 5. Solid line, \circ and \bullet as in Fig. 5; dotted and dashed lines (third-order polynomial fit to \circ and \bullet , respectively). (b) Expected analysis errors $\Delta\phi$. Solid line (theoretical); dotted line (first guess from data); and dashed line (second guess from data).

$15 < u_o < 20 \text{ m s}^{-1}$ where equally likely overestimated. The differences between first- and second-guess corrections are somewhat larger than for the mean wind speeds in Fig. 6(b). Large oscillations of the estimated analysis error again are outside the range in which the analysis renders results. Note that the correction for the bin width (35) is irrelevant here due to the large model error.

5. APPLICATIONS

As a further illustration, results of two case-studies are presented. In the first study random anemometer errors are assessed. It is shown that previous studies seriously overestimate the anemometer error owing to the finite bin width δu_o , and the implicit assumption of error-free observations. In the second study, systematic and random errors of ERS-1* scatterometer and SSM/I† wind speeds are estimated from collocations with buoy observations. The latter example illustrates difficulties in estimating observation errors.

* ESA (European Space Agency) Remote-sensing Satellite.

† Special Sensor Microwave Imager.

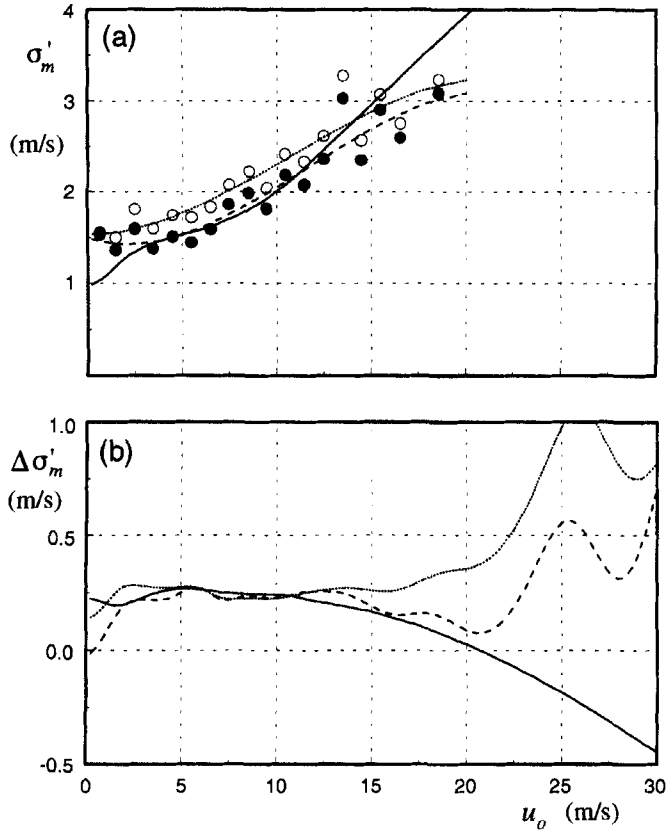


Figure 7. Like Fig. 6 but for the model error σ'_m and analysis error $\Delta\sigma'_m$.

(a) Anemometer errors

Anemometer instrument errors are usually estimated from a side-by-side comparison of identical instruments (Gilhausen 1987; Monaldo 1988). For the present study such anemometer data were provided by the National Data Buoy Center (NDBC) for five buoys for the period from December 1994 to February 1995 (Table 2). Because identical instruments are intercompared, the 'instrument' and 'model' error are identical, and the geometric mean regression should be used. Furthermore, biases β should be attributed to both anemometers ($\beta = \frac{1}{2}(\bar{u}_2 - \bar{u}_1)$, the suffices identify the two anemometers). The resulting biases, regression coefficients and random errors are presented in Table 2. For all buoys the biases are of the same order of magnitude as the random error or even larger. This implies that (for a single anemometer) calibration errors are of the same order of magnitude as the random error. In most studies results of many anemometers are used. Assuming that the calibration errors are uncorrelated, they then become part of the overall random instrument error (as in the last line of Table 2).

In his Fig. 3, Gilhausen (1987) estimates random anemometer errors as a function of the wind speed using a BA analysis. Applying a similar analysis to the present data set reproduces his results closely (\circ in Fig. 8). However, because the anemometers are very accurate, this analysis is contaminated by the bin width δu_o , in particular for low wind speeds (compare + and \circ in Fig. 8). Furthermore, the BA analysis implicitly assumes that

TABLE 2. BIASES β , REGRESSION COEFFICIENTS b_{gm} AND MEAN RANDOM ERRORS $\overline{\sigma}'_2 \approx \overline{\sigma}'_1$ FOR DUPLICATE ANEMOMETERS AT FIVE BUOYS

Buoy	Number of observations	β (m s^{-1})	b_{gm}	$\overline{\sigma}'_2$ (m s^{-1})
46035	1068	-0.10	0.944	0.20
46001	1333	0.20	1.046	0.11
51001	2147	0.10	1.027	0.09
42001	2132	-0.08	0.970	0.18
44008	249	-0.03	0.993	0.04
All	6929	0.03	0.992	0.24

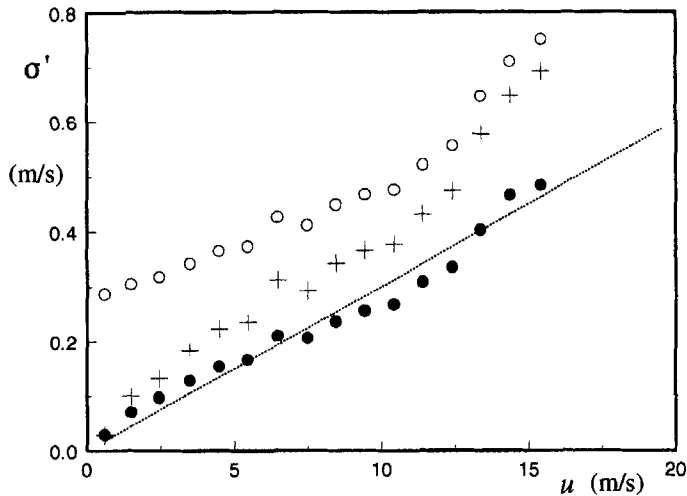


Figure 8. Estimates of random anemometer errors σ' as a function of the wind speed u for the anemometer data of Table 2. \circ (BA analysis); $+$ (BA analysis with bin-width correction); \bullet (ECBA analysis, iteratively applied until $\sigma'_1 = \sigma'_2$); and dotted line ($\sigma' = 0.03u$).

one anemometer is free of errors, and hence assigns all errors to the second anemometer. It is more reasonable to assume that both anemometers have identical errors ($\sigma'_1 = \sigma'_2$). Iteratively applying an ECBA analysis to satisfy this equality reduces the estimate of the random anemometer error by another 30% (compare \bullet and $+$ in Fig. 8), close to the expected reduction by a factor $\sqrt{0.5}$. The resulting instrument error is approximately 3% of the wind speed. Such a linear relation might be expected because the error is dominated by calibration differences of individual anemometers.

(b) Wind speeds from satellites

As a second illustration biases of wind speed retrievals from satellites are estimated. Considered are fast delivery ERS-1 scatterometer wind speeds (Offiler 1994) and SSM/I F13 wind speeds according to Goodberlet *et al.* (1990). These wind speed retrievals have been collocated with deep-ocean buoy observations for the period from December 1994 to February 1995 using a collocation radius of 50 km and 30 min. This resulted in 454 and 1202 collocations, respectively*.

* These collocations were performed as part of the validation study of a new wind-wave forecast system at the National Centers for Environmental Prediction. More details will be presented elsewhere.

TABLE 3. ESTIMATES OF BUOY OBSERVATION ERRORS IN PERCENT FOR SATELLITE WIND SPEED RETRIEVALS

	Collocation		Scale representativeness	Total relative error γ_0		
	space	time		low	best	high
ERS-1 scatterometer	5-6	3	3.5-6.5	7.7	8.8	9.9
SSM/I F13	8-11	3	5-8	10.5	12.3	14.3

The total relative error γ_0 is based on a global mean wind speed of 8 m s^{-1} , and represents a low, best and high estimate, respectively. γ_0 includes a 3.5% instrument and round-off error.

Observation errors for marine wind speeds have been investigated by, for instance, Brown (1983), Pierson (1983), Gilhausen (1987) and Monaldo (1988). Several types of observation errors can be distinguished; for instance (i) instrument errors, (ii) round-off errors due to data transmission and archiving, and (iii) mismatch errors in collocation and in representative scales (known as the representativeness error in data assimilation (Lorenz 1986)). An honest validation of a model or a retrieval algorithm considers observations which are representative for the validated parameter. Representativeness errors therefore should be considered as a part of the observation error.

The instrument error of the buoy observation is estimated in the previous subsection as 3% of the wind speed. The buoy data used here were archived with an accuracy of 0.5 m s^{-1} , which corresponds to an additional random observation error of approximately 0.15 m s^{-1} . The minimum observation $\sigma'_{o,\min}$ error of this buoy data is thus

$$\sigma'_{o,\min} = \sqrt{(0.03 u_o)^2 + (0.15 \text{ m s}^{-1})^2}, \quad (43)$$

which corresponds to an error of 3.5% for a global mean wind speed of 8 m s^{-1} . Representativeness errors arise due to collocation and scale mismatch errors in space or time. Estimates for such errors can be obtained from Pierson (1983), Gilhausen (1987) and Monaldo (1988). Estimates of collocation errors can be obtained directly from these papers, using calculated average collocation distances of 16 and 25 km for ERS-1 and SSM/I data, respectively. Scale representativeness errors are more difficult to estimate. A detailed discussion of such estimates will be presented elsewhere. For the present study it suffices to say estimates of observation errors almost always incorporate significant uncertainties.

A tally of error estimates is presented in Table 3. These mean errors are sufficient to correct the linear regression. The ECBA analysis, however, requires an estimate of the observation error as a function of the wind speed. Scale errors increase (approximately linearly) with wind speeds (Pierson 1983; Monaldo 1988). Unfortunately, collocation errors have not been assessed as a function of wind speed. It appears natural to assume that the overall error increases approximately linearly for higher wind speeds, and is finite for small wind speeds. This suggests a shape of the observation error similar to Eq. (43)

$$\sigma'_o = \sqrt{(\alpha \sigma'_o)^2 + (\gamma u_o)^2}, \quad (44)$$

where $0 < \alpha < 1$. Furthermore requiring that the average error fraction γ_0 is reproduced for the mean wind speed \bar{u}_o , the asymptotic error fraction γ becomes

$$\gamma = \gamma_0 \sqrt{1 - \alpha^2}. \quad (45)$$

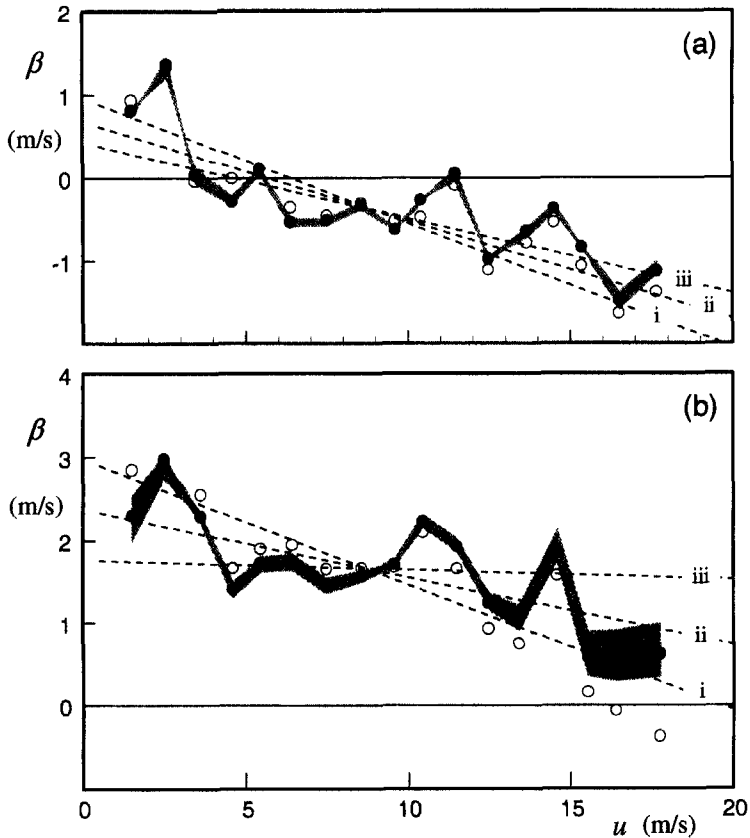


Figure 9. Estimates of biases β of wind speed retrievals from (a) ERS-1 fast delivery scatterometer and (b) SSM/I F13 (Goodberlet *et al.* 1990). Dashed lines (linear regression analysis (i) conventional, (ii) error-corrected conventional and (iii) geometric mean (symmetric slope); \circ (BA analysis); and \bullet (ECBA analysis with best estimate γ_0 from Table 3 and $\alpha = 0.7$ in Eq. (44)). Shaded area is the ECBA analysis for γ_0 ranging as in Table 3 and $0.6 < \alpha < 0.8$.

Somewhat arbitrarily, $\alpha = 0.7$ is assumed, corresponding to $\gamma/\gamma_0 = 0.71$. To assess uncertainties in this assumption, a range of $0.6 < \alpha < 0.8$ ($0.8 > \gamma/\gamma_0 > 0.6$) will be considered in the calculations.

Bias estimates for satellite retrieved wind speeds are presented in Fig. 9. For both satellites significant biases are found. Whereas these results are interesting in their own right, the differences between analysis techniques are more interesting in the present context. Regression line (ii) represents the best linear estimate of the biases. The conventional linear regression lines (i) and the geometric mean regression (iii) both deviate significantly from this best guess, in particular for the SSM/I wind speed retrievals. The conventional regression by definition underestimates the regression coefficient, and therefore overestimates the dependence of the bias β on the wind speed u in both cases. The geometric mean regression (iii) underestimates the regression coefficient for both data sets, indicating that the random error of both retrieval algorithms is larger than the corresponding observation error. In particular for the SSM/I wind speed retrievals, the geometric mean regression behaves poorly. It suggests that the bias is nearly independent of the wind speed, whereas the error-corrected regression shows a significant dependency.

For the ERS-1 data (Fig. 9(a)), effects of the error correction are relatively small, as was the case for the linear regression analyses. For the SSM/I data (Fig. 9(b)), the corrections are larger. For the latter data the BA analysis appears to suggest somewhat nonlinear bias behaviour for high wind speeds, which is less apparent in the ECBA analysis. The shaded areas in Fig. 9 represent the uncertainty in the ECBA analyses introduced by the uncertainty in the estimation of the observation errors. Whereas this uncertainty is appreciable, it does not influence the effects of the error correction qualitatively, nor does it appear to be relevant for the intercomparison of the separate analysis techniques. For the present data, the point-to-point variability of the results is larger than the error correction (● versus ○) and its uncertainty (shaded area). This variability represents the sampling error, which is significant for the small data sets considered here. For larger data sets the sampling error will become smaller, and in many cases the sampling error will become negligible compared with the analysis error (additional examples will be presented elsewhere).

6. DISCUSSION

The present study addresses effects of errors in observations on the results of validation studies, in particular where linear regression or BA analyses are used. Though it considers marine surface wind speeds its results are applicable to a wide range of validation studies in many fields of theoretical and applied research.

Linear regression analyses have been used for decades in validation studies, and many text books on the subject can be found. It has long been known that a conventional regression analysis (17) is valid only if the error in the observation is negligible. Unfortunately, effects of observation errors are rarely discussed in text books or in validation studies. The linear estimate of the functional behaviour of a modelled wind speed $U_m = \Phi(U_o)$ is bounded by the conventional regression (17), and the 'inverse' regression (20). In special cases, where the ratio of observation and model errors can be estimated, more advanced regression techniques can be used. Often used is the geometric mean or symmetric slope regression, which implicitly assumes that model and observation errors are similar. Whereas this regression can be used successfully for selected studies (see section 5(a)), this regression is also prone to errors if applied indiscriminately (see Table 1 or, for example, Lindley (1947)). If the observation error can be estimated, the conventional regression can be corrected as in Eq. (33) to obtain a best possible linear estimate of the functional relation.

A more advanced way to estimate the functional behaviour of a model is the BA analysis, which is intended to identify the functional behaviour and the random error of a model as a function of the wind speed. However, a BA analysis also incorporates systematic errors due to the convolution of the data with the observation error. Three major effects of observation errors are:

- (i) The BA analysis systematically underestimates extreme wind speeds;
- (ii) It systematically overestimates random model errors as observation errors are attributed to the model. If the model is accurate and the bin width large, the bin width also artificially increases the estimate of the random model error; and
- (iii) The analysis error is nonlinear, which can erroneously suggest or mask nonlinear model behaviour.

If the observation error can be estimated, the BA analysis can be corrected, as is demonstrated in sections 4 and 5.

Corrections of the regression and BA analyses stand or fall with a good estimate of the observation error. The observation error includes both the instrument error and a

representativeness error, as is discussed in section 5(b). This representativeness error is often much larger than the instrument error. For instance, Table 3 shows representativeness errors of two to four times the instrument error, and similar ratios can be found for other meteorological parameters (e.g. Kitchen 1989; Ingleby 1995). Ignoring this error can easily lead to the erroneous conclusion that a conventional analysis can be used because the instrument error (instead of the observation error) is much smaller than the model error. For practical studies the observation error will never be known exactly. It is therefore prudent to address the sensitivity of results to the assumed model error (as in Fig. 9). Furthermore, additional studies aimed at refining our knowledge of observation errors appear appropriate, both in the context of the present validation techniques and in the context of data assimilation.

If the observation error can only be estimated crudely, the question arises if the present techniques should be used. Establishing the impact of the uncertainty of the observation error as discussed above then becomes crucial. For the ECBA analysis, which is sensitive to both the overall error and its distribution, this might well mean that no definite conclusions can be drawn. This in itself would be an important finding, as it implies that conventional BA techniques also cannot be trusted. The error-corrected regression analysis is more 'robust', as it depends on a bulk error estimate only, and therefore generally will be more conclusive.

The present error-corrected regression technique is particularly useful if forecast systems are validated. In such systems, model errors typically grow with forecast time, whereas observation errors remain constant. The normalized regression coefficient b_{gm}/B_o of a geometric mean regression then systematically increases during the forecast (Eq. (23)). Such a systematically changing analysis error should obviously be avoided. The analysis errors of a conventional regression analysis are systematic and constant during the forecast (Eq. (19)). Such analysis errors are less detrimental in the validation of forecast systems. The present error-correction technique cannot remove the latter error completely if the observation error can only be estimated crudely. If the estimated error is defined as $(1 + \epsilon)\langle S'_{oo} \rangle$, where ϵ identifies the relative error in the estimation of the random observation error, the normalized error-corrected regression coefficient $b_{o,c}$ becomes

$$\frac{b_{o,c}}{B_o} = \left(1 - \frac{\epsilon \langle S'_{oo} \rangle}{S_{oo}} \right)^{-1}. \quad (46)$$

This implies a systematic error that does not change during the forecast, and is smaller than the error of the conventional regression as long as $\epsilon < 1$. Thus, the error-corrected regression is preferable for the validation of forecast systems as long as the observation error variance is not overestimated by more than a factor of 2.

An interesting property of the ECBA analysis is that it can retrieve random model errors σ'_m which are smaller than the observation error σ'_o . This suggests that accurate validation results can be obtained from poor quality observations. However, two potential problems occur in this scenario. First, large observation errors imply a large correction of the ECBA analysis relative to the BA analysis, making effects of the uncertainty in the observation error potentially sizeable. Second, the error corrections suggested in section 3 are based on a convolution of the *observed* distribution with an error pdf. This implies that the retrieved random model error corresponding to the estimated moments $\hat{m}_0, \hat{m}_1, \dots$, etc. is always positive, even if the observation error is grossly overestimated and provides more variability than supported by the data. This has two consequences: (i) if the model error is significantly smaller than the observation error, the ECBA method is expected to overestimate the random model error, and (ii) the validity of the estimated observation error

should be tested independently, to assure that the data support this observation error. The mean model error $\bar{\sigma}'_m$ of the corrected linear regression (Eqs. (33) and (25)) provides such a test. This error becomes 0 if all variability in the data is explained by the observation error, and becomes undefined if the data cannot support the assumed observation error.

7. CONCLUSIONS

Conventional linear regression and bin-averaged validation techniques introduce systematic analysis errors if observation errors are not negligible. A conventional regression analysis then underestimates the regression coefficient and overestimates the random model error. More advanced regression techniques like the geometric mean (or symmetric slope) regression then also incorporate significant errors for many applications. A BA analysis then underestimates extreme (high) wind speeds, incorporates spurious nonlinearity, and overestimates random model errors. Example calculations with synthetic and real data suggest that such errors are generally not negligible in detailed validation studies of marine wind speeds. If the observation error as a function of the wind speed can be estimated, it is possible to remove the systematic errors from the above validation techniques. Present knowledge of observation error is sufficient to apply the present error-correction techniques in many cases, but our understanding of observation error could be improved significantly. Finally, it is argued that the geometric mean (symmetric slope) regression should not be used to validate forecast systems, because its analysis errors are expected to be functions of the forecast time.

Although this paper explicitly deals with wind speeds, its results are expected to be valid for a wide range of validation studies in many fields of research.

ACKNOWLEDGEMENTS

The author acknowledges with pleasure that the idea for the cut-off normal distribution (15) was supplied by R. J. Purser. The author thanks E. Meindl and R. Strahan of NDBC for supplying the anemometer data of Table 2 and Fig. 8, and W. G. Gemmill, V. Krasnopolsky, R. J. Purser, W. J. Pierson, N. B. Ingleby and an anonymous referee for constructive comments on early versions of this manuscript. This research was undertaken under the auspices of the UCAR Visiting Scientists Program at the Ocean Modelling Branch (OMB) of the NCEP. The manuscript is designated OMB Nr 119.

REFERENCES

- | | | |
|--|------|---|
| Abramowitz, M. and Stegun, I. A. | 1973 | <i>Handbook of mathematical functions</i> , 9th print. Dover Publications, New York |
| Brown, R. A. | 1983 | On a satellite scatterometer as an anemometer. <i>J. Geophys. Res.</i> , 88 , 1663–1673 |
| Berkson, M. D. | 1950 | Are there two regressions? <i>J. Am. Statist. Ass.</i> , 45 , 164–180 |
| Draper, N. R. and Smith, H. | 1981 | <i>Applied regression analysis</i> . Wiley |
| Gilhausen, D. B. | 1987 | A field evaluation of NDBC moored buoy winds. <i>J. Atmos. Ocean. Technol.</i> , 4 , 94–104 |
| Goodberlet, M. A., Swift, C. T. and Wilkerson, J. C. | 1990 | Ocean surface wind speed measurements of the Special Sensor Microwave/Imager(SSM/I). <i>IEEE Trans. Geoscience and Remote Sensing</i> , 28 , 823–827 |
| Hinton, B. B. and Wylie, D. P. | 1985 | A correction for the errors in ship reports of light winds. <i>J. Atmos. Ocean. Technol.</i> , 2 , 353–356 |
| Ingleby, N. B. | 1995 | Assimilation of station level pressure and errors in station height. <i>Weather and Forecasting</i> , 10 , 172–182 |
| Jolliffe, I. T. | 1990 | Principal component analysis: a beginner's guide. I. Introduction and application. <i>Weather</i> , 45 , 375–382 |

- Kitchen, M. 1989 Representativeness errors for radiosonde observations. *Q. J. R. Meteorol. Soc.*, **115**, 673–700
- Lindley, D. V. 1947 Regression lines and the linear functional relation. *J. R. Statist. Soc.*, ser. B, **9**, 218–244
- Lorenc, A. C. 1986 Analysis methods for numerical weather prediction. *Q. J. R. Meteorol. Soc.*, **112**, 1177–1194
- Monaldo, F. 1988 Expected differences between buoy and radar altimeter estimates of wind speed and significant wave height and their implications on buoy–altimeter comparisons. *J. Geophys. Res.*, **93**, 2285–2302
- Offiler, D. 1994 The calibration of ERS–1 satellite scatterometer winds. *J. Atmos. Ocean. Technol.*, **11**, 1002–1017
- Pierson, W. J. 1983 The measurement of synoptic scale wind over the ocean. *J. Geophys. Res.*, **88**, 1683–1708
- Ricker, W. E. 1973 Linear regression in fishery research. *J. Fishery Research Board of Canada*, **30**, 409–434