

STATISTICAL EVALUATION OF NEWLY ACQUIRED ENVIRONMENTAL CHEMICAL ANALYSIS DATA AS A VALIDATION TOOL FOR THE URANIUM MILL TAILINGS REMEDIAL ACTION (UMTRA) PROJECT

Kimberly C. Smith and Paul E. Zietz

Roy F. Weston, Inc.

2155 Louisiana NE

Suite 10,000

Albuquerque, NM 87110

Mark R. Minter

SHB AGRA, Inc.

2155 Louisiana NE

Suite 10,000

Albuquerque, NM 87110

ABSTRACT

A statistical comparison of new data with existing data is being used as an effective tool to augment standard data validation practices on the Department of Energy Uranium Mill Tailings Remedial Action (UMTRA) Project. The automated system allows rapid and reliable detection of anomalies in newly acquired chemical analysis data. Three separate statistical methods have been combined to meet the demanding requirements of environmental chemical data, including age-weighted trend analysis, downweighting of historical anomalies, and iterative recalculation of substitute values for below-detection data.

INTRODUCTION

The U.S. Department of Energy Uranium Mill Tailings Remedial Action (UMTRA) Project is responsible for the assessment and remediation of health and environmental hazards associated with uranium processing activities at 24 sites in 10 states. An important component of the Project is monitoring of groundwater and surface water quality on and around the sites. Since the early 1980s, 42 inorganic and 7 radiochemical constituents have been monitored at several sampling locations on each site. The UMTRA water-quality database currently contains over half a million measurements, and approximately 50,000 new measurements will be acquired this year.

Standard environmental chemical analysis data validation practices include verifying unbroken chain of sample custody, confirming adherence to holding-time requirements, and evaluating quality-control sample results. Unfortunately, many serious problems such as dilution factor errors, transcription errors, sample or process contamination, and incorrect reporting units cannot be identified in this way. To address these problems, the UMTRA Project introduced during 1992 a statistical comparison of new data with historical data as an additional tool in the data validation process. The statistical comparison software was run during the initial electronic loading of the new data into the project database. Newly acquired data that departed in a significant way from historical measurements were identified in a "suspected anomaly report".

The anomaly detection program was successful in identifying many analytical and reporting errors that otherwise would have gone undetected. Because the statistical comparison program was automated, suspected anomalies were identified quickly, allowing for reanalysis of questionable results in many cases. However, the stationary, control-chart approach used in the anomaly detection program was inadequate when applied to historical data sets that exhibited trends in concentration or activity levels over time, had censored ("less than" or "nondetect") data, or contained anomalies among the older data. Therefore, a new statistical

component has been developed for use on the UMTRA project. The software incorporates these special features:

- The algorithm allows for trends over time in concentration or activity, such that an expected increase or decrease in a parameter value does not result in the signal of an anomaly. This feature also identifies data that do not conform to established trends so that these data can be given special scrutiny.
- Data acquired early in the project were not subjected to modern validation procedures, and are weighted so that for purposes of anomaly detection, these old data have less influence than more recent data.
- The influence of historical outliers in the database is restricted. This is especially important when very old data exist and validation issues relating to such outliers cannot be resolved satisfactorily.
- The software incorporates methods specifically designed to handle less-than-detection-limit data.

UMTRA's anomaly detection procedures are described in this paper, and several examples illustrate the advantages of the new statistical features. The software is currently in a testing phase and is expected to be fully operational within the next year.

DATA INPUT

Newly acquired data are electronically loaded into the UMTRA Project database. This saves a great deal of time and eliminates one activity that can introduce transcription errors. The batch-loading software automatically imports all new data from a sampling round into the anomaly detection computer program, and prints a "suspected anomalies report" when the statistical comparisons with historical data are complete.

STATISTICAL METHODS

The historical database for a parameter at a sampling location consists of ordered pairs (x,y) of information from all previous sampling rounds.

- x: date when water sample was collected, in decimal notation. For example, a sampling date of April 15, 1990 corresponds to $x = 90.288$.
- y: parameter value, typically reported in milligrams per liter (mg/L) for inorganic constituents, and picocuries per liter (pCi/L) for radiochemical parameters.

If an inorganic parameter value falls below the contract required detection limit (CRDL), the ordered pair will appear as $(x, < y)$.

In this paper, the symbol n represents the number of previous measurements available for a parameter at a sampling location. The analytical results are ordered by date of sampling round, so that (x_1, y_1) is the earliest and (x_n, y_n) the most recent data in the historical database. These historical data are used to compute an "acceptance interval" for a newly acquired analytical result. If the new measurement, denoted here by (x_{new}, y_{new}) , does not lie within the acceptance interval, the result is identified as a suspected anomaly and is subjected to further scrutiny by the technical staff. Calculation of the acceptance interval proceeds in a number of steps, described below.

Age-Based Weighting

More recent measurements are inherently better predictors of future measurements than are data collected several years in the past. Therefore, data are weighted based on their age relative to the most recent sampling round, that is on $x_n - x$. The form of weight function $w(x)$ being employed at UMTRA is the s-shape logistic function

$$w(x) = 1/\{1 + c \cdot (x_n - x)^p\} \quad (1)$$

where $c > 0$ and $p > 0$ are constants that influence the shape of the function. The logistic function assigns the most recent sampling period a weight of one, and previous sampling dates progressively less weight. Other age-weight functions are possible, but the logistic function was selected for use on the UMTRA Project because of the wide range of shapes that the logistic function can achieve by choice of the shape constants c and p , including the possibility of equal weighting of all data ($c = 0$).

A practical method of selecting values for constants c and p involves first selecting ages, $A_{0.5}$ and $A_{0.1}$ for which one wishes to assign weights 0.5 and 0.1 respectively. Then constants c and p are calculated as follows:

$$p = \log(9)/\{\log(A_{0.1}) - \log(A_{0.5})\} \quad (2)$$

$$c = (1/A_{0.5})^p \quad (3)$$

For example, we have determined that for the UMTRA anomaly-detection procedure it is reasonable to assign weight 0.5 to data that are 4 years older and weight 0.1 to data that are 10 years older than the most recent sampling round. Then $A_{0.5} = 4$ and $A_{0.1} = 10$. With these choices,

$$p = \log(9)/\{\log(10) - \log(4)\} = 2.4 \quad (4)$$

$$c = (1/4)^{2.4} = 0.036 \quad (5)$$

$$w(x) = 1/\{1 + 0.036 \cdot (x_n - x)^{2.4}\} \quad (6)$$

The UMTRA weight function is shown on Fig. 1 along with some other choices of c and p .

Trend Estimation Using Age-Based Weighting

Changes in concentration or activity levels of a parameter over time are common at UMTRA sites for which groundwater contamination plumes are present. Linear trends and a

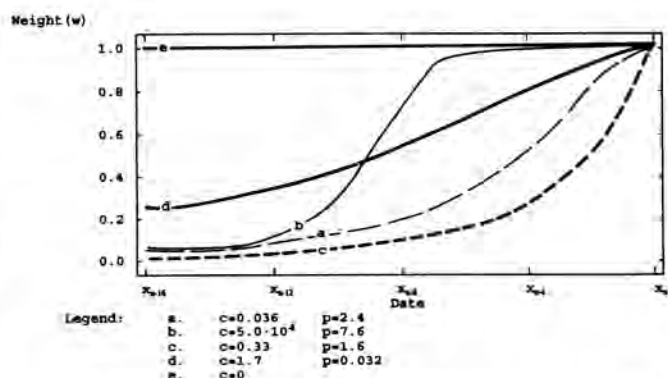


Fig. 1. Logistic age-weight function $w(x) = 1/\{1 + c(x - x_n)^p\}$

variety of nonlinear trends have been observed in UMTRA data sets. We have elected to estimate time trends by fitting a straight line through the age-weighted data, rather than to try to fit curves. Most curves are well approximated by a line, at least within a narrow range. The use of age-weighting in linear regression produces a straight line that mimics the linear behavior of the more recent data. The resulting fit through the older data may be poor, but the method is well suited to our task of constructing an acceptance interval for a new measurement.

The model for a straight line is $Y = \alpha + \beta \cdot x + \sigma \cdot \text{error}$, where α and β are the y-intercept and slope of the line respectively, σ represents the average amount by which real data will deviate from a perfect straight line, and the error is a random observation from some standard probability distribution F . It is common, but not necessary, to assume that this probability distribution is normal.

Model constants α , β , and σ are estimated from the available data so that the resulting line "fits" the data as closely as possible. If (a, b, s) are candidate estimates for (α, β, σ) , then discrepancies between the actual measurements and the candidate line $Y = a + b \cdot x$ are called "residuals". The i^{th} residual, r_i , is the distance that the data point (x_i, y_i) lies above or below the line on an x-y plot, $r_i = y_i - (a + b \cdot x_i)$. Candidates (a, b, s) are good estimates of (α, β, σ) to the extent that the residuals are small in some overall sense. A mathematical formula that describes how to compute the overall lack of fit between data and a fitted model is called a loss function. Upon selection of a loss function, the "best" parameter estimates (a, b, s) are those that produce minimum loss. For example, the well-known "ordinary least squares" (OLS) procedure computes overall loss by simple addition of the squared residuals:

$$\text{LOSSOLS} = \sum_i r_i^2 \quad (7)$$

In order to incorporate age-weighting, the OLS loss function is altered to

$$\text{LOSSAGE WEIGHTED} = \sum_i w(x_i) \cdot r_i^2 \quad (8)$$

The greater the weight assigned to a data pair (x, y) , the greater is the penalty assessed to a discrepancy between fitted line and measurement. Consequently, the weighted loss function results in a line that favors the more heavily weighted data points, that is, the more recently acquired data.

For the UMTRA Project, we have made two modifications to the age-weighted loss function. The first modification restricts the influence of "outliers" in a data set. The second modification is necessary to incorporate "less than detection limit" data. These modifications are described in the following two sections. Many of the technical details of these

modifications have been omitted here for the sake of clarity. Interested readers can find additional information in the references.

Downweighting of Outliers

An outlier is a measurement that clearly lies away from the trend of the other measurements in a data set. Outliers may be erroneous measurements, or may represent true but fleeting excursions in water quality. In either case, a highly unusual value will pull the fitted line towards itself and away from the more reasonable values. The use of age-weighting provides a great deal of protection against outliers that occur in older data, since older data values are assigned smaller weights anyway. However, in spite of improved quality assurance practices on the UMTRA Project in recent years, occasional outliers still occur in our more recently acquired, and more heavily weighted, data.

In an effort to make trend analysis more "robust", that is, resistant to the influence of outliers, an additional weighting factor is incorporated into the loss function. If, after calculating the best age-weighted line through the data, any measurements lie an excessive distance from that line, these measurements are downweighted, then the age-weighted fitting procedure is repeated. The process of fitting, then downweighting outliers, may be repeated several times.

The UMTRA algorithm defines an outlier as any age-weighted measurement that lies more than ± 1.5 standard deviations from the fitted line. The outlier-weighting factor applied to an outlier is proportional to the observed excess over ± 1.5 standard deviations. For example, if an age-weighted measurement is 1.8 standard deviations above or below the fitted line, its outlier-weight for the next iteration is $(1.5/1.8) = 0.83$. The outlier-weight is multiplied by the age-weight for the measurement. Ideally, the allowable number of standard deviations should increase with sample size, but preliminary use of 1.5 standard deviations has proved satisfactory on UMTRA data sets with 5 to 15 measurements.

Incorporation of Censored Values Through Iteration

A simple but unsophisticated approach for censored data is to assign the measured concentration some arbitrary value between zero and the detection limit itself, and proceed as if that arbitrary value is the true measurement. For example, the EPA recommends substitution of one-half the detection limit (DL/2) for each censored value, provided there are relatively few censored values in the data set (1). This recommendation is best suited to cases where the same detection limit is applied to each sampling round. However, it is common in UMTRA Project data sets for there to be two to five different detection limits in a particular data set. Improvements in analytical procedures have often resulted in a reduction of the detection limit over time. Conversely, plume movement can result in new matrix interferences and correspondingly increased detection limits due to required dilutions. When detection limits show time trends, the use of a substitute value derived from the detection limit can introduce a bias into the line-fitting process by creating the appearance that the concentration itself is decreasing or increasing with time. Finally, substitution methods become increasingly untenable as the percentage of censored values increases.

A data pair ($x, < y$) indicates that the true (but unobservable) measurement of concentration of the parameter on date x is equal to some value between 0 and y . Under the assumed

linear model, the unobservable measurement is $\alpha + \beta \cdot x + \sigma \cdot \text{error}$, constrained as follows:

$$0 \leq \alpha + \beta \cdot x + \sigma \cdot \text{error} < y \quad (9)$$

or, upon rearrangement,

$$-(\alpha + \beta \cdot x)/\sigma \leq \text{error} < (y - \alpha - \beta \cdot x)/\sigma \quad (10)$$

The UMTRA anomaly detection algorithm uses the two relationships above and information provided by all the data in order to select an appropriate substitute value for a censored measurement. The procedure is iterative, as outlined below.

1. Get the first set of minimum loss estimates (a_1, b_1, s_1) by substituting DL/2 for each censored measurement.
2. The expected value of an error, given that the error is constrained to lie between $-(a_1 + b_1 \cdot x)/s_1$ and $(y - a_1 - b_1 \cdot x)/s_1$ equals

$$\frac{\int_{-(a_1 + b_1 \cdot x)/s_1}^{(y - a_1 - b_1 \cdot x)/s_1} z \cdot dF(z)}{\int_{-(a_1 + b_1 \cdot x)/s_1}^{(y - a_1 - b_1 \cdot x)/s_1} dF(z)} \quad (11)$$

where $F(z)$ is theoretical probability distribution responsible for the random deviations of data from a straight line. Under the assumption that $F(z)$ is the normal distribution, this expected error, which we'll call error_1 equals

$$\frac{\phi[(y - a_1 - b_1 \cdot x)/s_1] - \phi[-(a_1 + b_1 \cdot x)/s_1]}{\Phi[(y - a_1 - b_1 \cdot x)/s_1] - \Phi[-(a_1 + b_1 \cdot x)/s_1]} \quad (12)$$

where $\phi[\cdot]$ and $\Phi[\cdot]$ are the normal density and cumulative distribution functions respectively.

3. Calculate improved estimates of the true measurements associated with each censored value $< y$:

$$a_1 + b_1 \cdot x + s_1 \cdot \text{error}_1 \quad (13)$$
4. Substitute the improved estimates for the nondetects and calculate a second set of minimum loss estimates (a_2, b_2, s_2).
5. Repeat steps 2, 3, and 4 until successive estimates become stable.

The performance of this iterative re-estimation procedure for nondetects, based on an assumption of normally distributed errors, has been demonstrated to be relatively unaffected if the actual error distribution is nonnormal (2).

Computation of the Acceptance Interval

The minimum loss estimates (a, b, s) are used to predict a future concentration at time x_{new} :

$$y_{\text{predicted}} = a + b \cdot x_{\text{new}} \quad (14)$$

The acceptance interval for the future concentration is a range of acceptable values, centered on this prediction, whose width depends on the random variation in the previous measurements (s) as well as the number of previous measurements (n) used to form the prediction

$$\text{Acceptance interval} = y_{\text{predicted}} \pm s \cdot k_n \quad (15)$$

where k_n is a constant that changes with n .

Standard statistical methods choose k_n in such a way as to control the probability of false positives at an acceptably

low level. In this context, a false positive result occurs if a perfectly valid measurement y_{new} falls outside the interval as a result of random variation alone, and is therefore flagged as a potential anomaly. Valuable time is spent investigating the quality of these measurements. The strategy used to reduce the probability of false positives is to widen the acceptance interval by choosing a relatively large value for k_n . The downside of this strategy is that a wide interval increases the probability of a false negative result. A false negative result occurs when a truly anomalous measurement falls in the acceptance interval, and consequently is not flagged for further investigation. Since the specific purpose of the procedure is to identify potential anomalies, an excessively wide acceptance interval is counterproductive.

The k_n used in acceptance intervals for UMTRA Project data validation are chosen to keep the probability of a false negative (as opposed to a false positive) at an acceptably low level. The k_n shown in Table I allow for the identification of approximately 99 percent of all anomalous measurements that deviate from the underlying linear model by 4 or more standard deviations. That is,

IF:

$$|y_{\text{new}} - (\alpha + \beta \cdot x_{\text{new}})| \geq 4 \cdot \sigma$$

THEN:

$$\text{Prob}[y_{\text{new}} \text{ will fall outside acceptance interval}] \geq 0.99 \quad (16)$$

where the probability assumes normally distributed errors. For data sets where the error distribution is nonnormal, the false negative protection rate may be slightly greater or less than the nominal 99 percent. A high level of false negative protection requires a rather narrow acceptance interval, with the unavoidable tradeoff of a high false positive rate. However, Table I shows that the false positive rate decreases as the sample size increases, so that over time the same level of false negative protection can be achieved at progressively less cost.

The k_n values in Table I were calculated using the UMTRA age-weight constants ($c=0.036$ and $p=2.4$ mentioned above) under the assumption that sampling occurs annually. At UMTRA sites, sampling is not always done annually, so in practice the critical values k_n are calculated individually for each data set.

OUTPUT

A "suspected anomalies report" is generated after completion of the data comparison steps described above. This report provides the reviewer with the background information necessary to decide if the flagged data are true anomalies, and hence if reanalysis should be requested. Some of the information included in the report is listed below.

- The new datum, any data qualifiers associated with the new measurements, the reported detection limit and, for radiochemical parameters, the associated analytical uncertainty.
- The number of previously acquired data points for that parameter and sampling location and the percentage of nondetects in the historical data set.
- The historical maximum and minimum values for that parameter and sampling location.
- The last three values acquired, with acquisition dates, data qualifiers, reported detection limits, and analytical uncertainties, if applicable.

TABLE I
Critical Values k_n for Acceptance Intervals,
 $Y_{\text{predicted}} \pm s \cdot k_n$

Sample Size n	Critical Value k_n^a	Probability of False Positive
4	0.95	0.61
6	1.43	0.37
8	1.65	0.26
10	1.83	0.19
12	1.92	0.16
15	2.15	0.10
20	2.29	0.07

^aCritical values are based on Normal error distribution, age-weight parameters $c=0.036$, $p=2.4$, and annual sampling. These critical values will maintain false negative probability rate at approximately 0.01.

Once the new UMTRA anomaly procedure is fully tested and operational, the suspected anomalies report will be modified to include the calculated acceptance interval for the new datum. The final software will also generate graphs similar to those presented in the example section that follows.

EXAMPLES

Several examples are presented below to demonstrate the performance of the UMTRA anomaly detection algorithm. The examples are similar to water quality data sets from UMTRA sites, but sampling dates have been set at annual increments for simplicity. The acceptance intervals are calculated for results from a hypothetical sampling round on January 1, 1994, so $x_{\text{new}} = 94.0$.

Example 1

Figure 2 represents an actual set of sulfate measurements from a well directly downgradient of a former uranium processing site. Concentrations of sulfate have been decreasing steadily since sampling began in 1982. When data sets have no gross outliers, no censored values, and no obvious curvature to the time trend, the UMTRA fitting algorithm is essentially indistinguishable from ordinary least squares regression.

Figure 2 shows the fitted line, $y = 6434 - (59.98)(x)$, through the data and the predicted concentration for the hypothetical January 1, 1994 sampling round, $y_{\text{predicted}} = 6434 - (59.98)(94.0) = 796 \text{ mg/L}$. Since there are 12 previous annual sampling rounds, the critical value $k_{12} = 1.92$ from Table I and the estimated standard deviation around the fitted line, $s = 45$, are used for computation of the acceptance

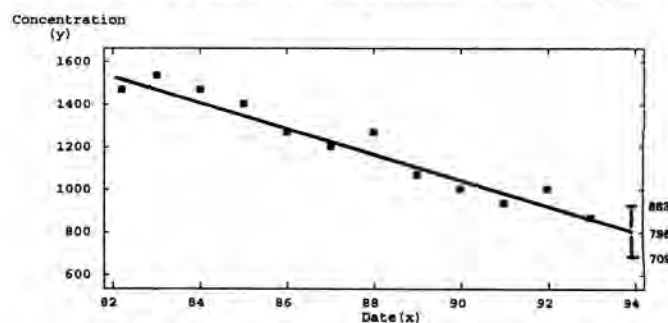


Fig. 2. Time-concentration plot for example 1.

interval: $796 \pm (1.92)(45) = 796 \pm 87$ or 709 to 883 mg/L. Therefore, if the sulfate measurement from the January 1, 1994 sampling round is less than 709 mg/L or greater than 883 mg/L, that result will be flagged as a suspected anomaly.

Example 2

This example demonstrates one benefit of age-based weighting. The linear trend in concentration shown in Fig. 3 is disrupted by an unusually high measurement that occurred in 1986, seven years earlier than the most recent sampling round of 1993. The UMTRA age-weighting constants $c = 0.036$, $p = 2.4$ downweight the measurement to 0.21. Once downweighted, the measurement lies within ± 1.5 standard deviations of the fitted line, so that additional outlier downweighting is not required. The age-weighted regression line and acceptance interval are essentially unaffected by the presence of the historical anomaly. Without the protection against historical anomalies provided by age-weighting, an acceptance interval for a future measurement may be excessively wide, increasing the chance that an anomaly in the newly acquired datum will go undetected. Figure 3 shows how much tighter the acceptance interval for the hypothetical January 1, 1994 datum becomes when age-weighting is used, compared to when it is not used.

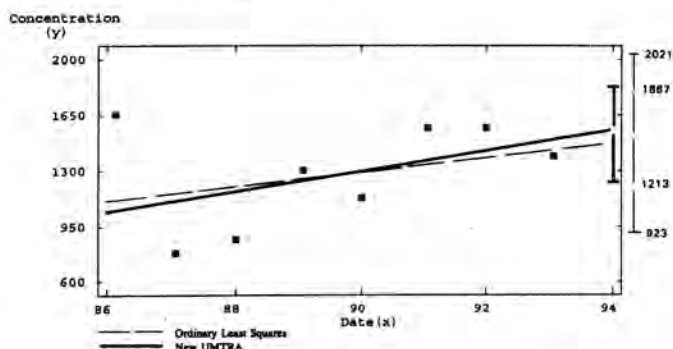


Fig. 3. Time-concentration plot for example 2.

Example 3

A less obvious benefit of age-weighting is demonstrated in this example. The concentration of the parameter depicted in Fig. 4 is neither increasing nor decreasing over the sampling periods shown. Improvements in analytical procedures, however, have increased the precision of measurements of the more recent sampling rounds. In addition to allowing recently acquired data more influence on the prediction of a future value, age-weighting also allows those data more influence on the measure of variability (s) around the line. As a result, the width of the acceptance interval tends to reflect the variation observed in recent measurements. Figure 4 compares the acceptance intervals produced with and without the age-weighting feature. The non-age-weighted interval seems too wide to reliably detect anomalies, considering the improvements in analytical precision of recent years.

Example 4

This example demonstrates the performance of the UMTRA anomaly-detection algorithm when applied to a curvilinear trend in concentration over time. The plot in Fig. 5 shows water quality measurements increasing over time, then leveling off, and beginning a slow decrease. A pattern like this may be associated with the passage of a discrete contam-

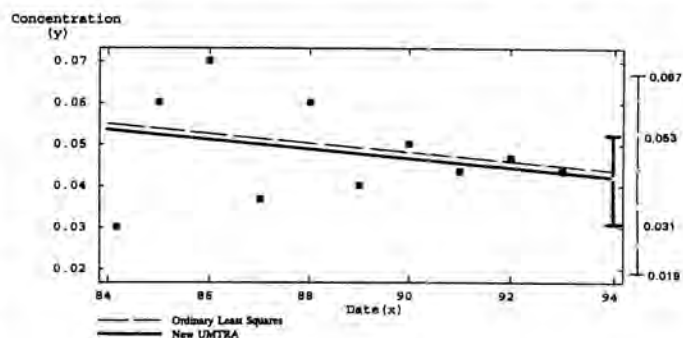


Fig. 4. Time-concentration plot for example 3.

ination plume through an aquifer. After age-weighted regression, the 1986 measurement lies 1.56 standard deviations from the fitted line. The outlier protection procedure becomes operative at this point, further downweighting the measurement by a factor of $(1.50/1.56) = 0.96$. Three additional iterations of fitting/downweighting produce the fitted line and acceptance interval shown on Fig. 6. The ordinary least squares line through the data produces a higher predicted value and a wider acceptance interval for the hypothetical January 1, 1994 measurement, due to the fact that ordinary least squares assigns all data equal weight for predicting the new measurement.

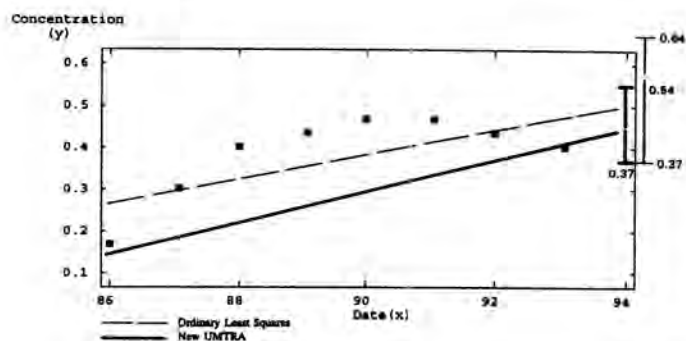


Fig. 5. Time-concentration plot for example 4.

Example 5

The data set presented in Fig. 6 contains two measurements that were reported as < 0.10 . These measurements are plotted on the graph with the symbol " $<$ " so that it is clear to the viewer that the unquantified measurement is equal to some value below the " $<$ ", but above the x-axis. The DL/2 rule produces substitute values of 0.05 for both measurements. The substitute values appear too low when compared to the other available measurements. The UMTRA fitting algorithm iteratively makes adjustments to the DL/2 substitution so that the substituted values become progressively more consistent with both the trend and the variability observed in the other data. The final fitted line and acceptance interval for the hypothetical new measurement on January 1, 1994, based on final substitute values of 0.78 and 0.86 for the two censored values, are indicated on Fig. 6.

Example 6

Example 6 shows another case for which the DL/2 rule is inadequate for analyzing trend in data sets with censored

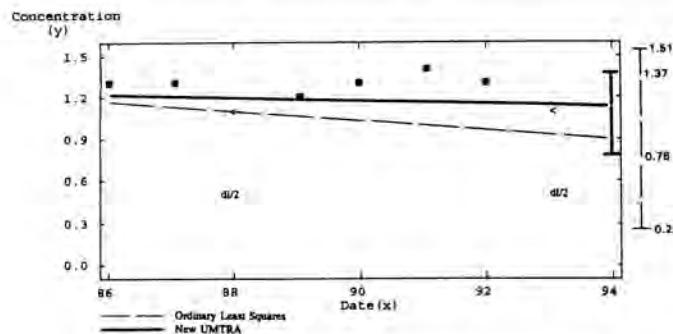


Fig. 6. Time-concentration plot for example 5.

measurements. The data in Fig. 7 are actual arsenic measurements from a well at an UMTRA processing site. A careful examination of the time-concentration plot reveals an increasing trend in the detection limit over time, but not necessarily the arsenic concentration itself. The DL/2 rule does not eliminate the appearance of trend in this example, and use of the DL/2 rule results in a higher predicted value for the hypothetical January 1, 1994 sampling round than was observed in any previous sampling round. On the other hand, the UMTRA anomaly procedure determines through iteration that the inflated detection limits of more recent times provide scant information on which to base a prediction. The prediction interval for the new January 1, 1994 arsenic measurement, 0.010 to 0.026 mg/L, essentially brackets the above-detection measurements in the data set. The newly acquired measurement would be deemed acceptable if it was either a) above detection and fell within the acceptance interval or b) was reported as below a detection limit greater than the lower limit of the acceptance interval, 0.010 mg/L.

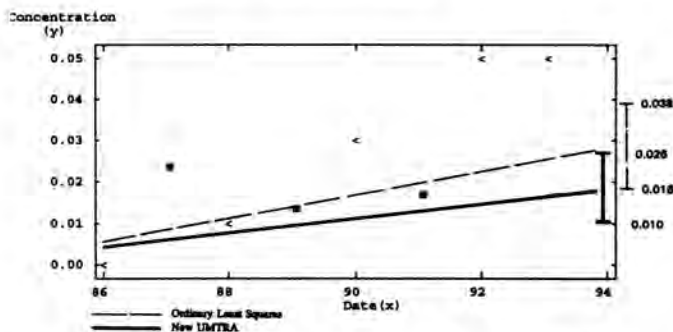


Fig. 7. Time-concentration plot for example 6.

Example 7

This final example demonstrates the value of the outlier down-weighting feature of the anomaly-detection algorithm. The hypothetical data set depicted in Fig. 8 includes a very unusual measurement for the most recent sampling round. By virtue of being the most recently acquired value in this particular data set, the datum receives an age-weight of one. In addition, endpoints inherently have greater leverage than other points in determination of a fitted line. As a consequence of these two factors, outliers in recent data have considerable influence on the prediction and acceptance interval for the next measurement.

The outlier protection strategy used in the UMTRA anomaly detection procedure substantially reduces the impact of recently acquired anomalies in data. The program

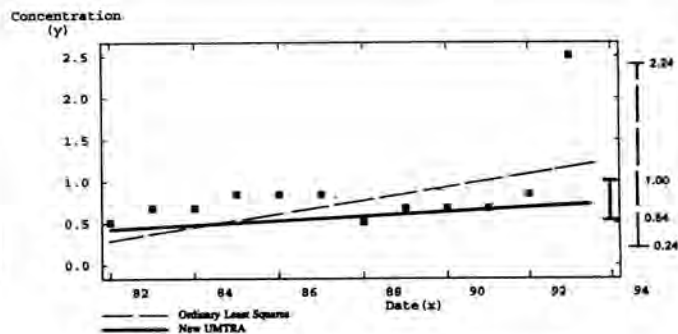


Fig. 8. Time-concentration plot for example 7.

required 20 iterations to reach a final solution, at which point the outlier-weight assigned to the anomaly was 0.32. Figure 8 shows the acceptance interval for the hypothetical measurement on January 1, 1994. Ordinary least squares regression on the data again fails to produce a useful acceptance interval for detection of anomalies in newly acquired data.

The example underscores the importance of a statistical anomaly detection program. A well-designed algorithm that identifies anomalies as they occur in newly acquired data will, over time, significantly reduce the number of problematic outliers in the database. Future users of the data will be spared the burden of making independent evaluations about the quality and useability of such measurements.

CONCLUSION

Given the many regulatory requirements and possible risks to human health and the environment associated with environmental restoration projects, early detection of anomalous results has social as well as scientific value by reducing the potential for issuance of spurious data to regulatory agencies and the public. A computerized anomaly detection procedure can be a powerful tool in data validation where large quantities of data are collected over extended periods of time. One or more of the components of the UMTRA program may be useful to other organizations that engage in long-term environmental monitoring or remediation activities. Age-weighted trend analysis is both easy to understand and easily automated, as no iteration is required to produce estimates. Additional outlier protection is also beneficial if older, less reliable data will be used for validation of newly acquired data. Projects that accumulate large quantities of censored data should also consider methods other than simple substitution for using "less than" data in data validation efforts. The censored regression component of the UMTRA procedure, although conceptually, mathematically, and computationally the most difficult feature to implement, is proving invaluable to the UMTRA data validation effort.

The several features incorporated into the UMTRA program are, individually, well-known and exhaustively tested statistical methods (3-6). However, the application of these methods to data validation, the simultaneous use of several methods on the same data set, the automation of the methods for use on literally thousands of data sets without user interface, and the emphasis on false negative protection in acceptance interval construction, together make this particular program unique. Consequently, the program is undergoing rigorous testing on actual and simulated data to identify computational limitations, and to verify that individual results

generally agree with the best professional judgment of the UMTRA staff who will use and interpret program output.

REFERENCES

1. U.S. ENVIRONMENTAL PROTECTION AGENCY, "Statistical Analysis of Ground-Water Monitoring Data at RCRA Facilities-Interim Final Guidance," (February, 1989).
2. G. HELLER and J.S. SIMONOFF, "A Comparison of Estimators for Regression with a Censored Response Variable," Biometrika 77(3): 515-520 (1990).
3. J. NETER, W. WASSERMAN, and M. H. KUTNER, Applied Linear Statistical Models: Regression, Analysis of Variance, and Experimental Designs, 2nd Edition, Irwin, Homewood, IL, (1985).
4. P.J. HUBER, Robust Statistics, New York, Wiley (1981).
5. J. SCHMEE and G.J. HAHN, "A Simple Method for Regression Analysis with Censored Data," Technometrics 21(4): 417-432 (1979).
6. M. AITKIN, "A Note on the Regression Analysis of Censored Data," Technometrics 23(2): 161-163 (1981).