

CONF-8510151--6

DE86 001335

DETECTING ISOTOPIC RATIO OUTLIERS

Charles K. Bayne

Computing and Telecommunications Division
Oak Ridge National Laboratory*
Oak Ridge, Tennessee 37831

David H. Smith

Analytical Chemistry Division
Oak Ridge National Laboratory*
Oak Ridge, Tennessee 37831

MASTER

By acceptance of this article, the publisher or recipient acknowledges the U.S. Government's right to retain a nonexclusive, royalty-free license in and to any copyright covering the article.

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

*Operated by Martin Marietta Energy Systems, Inc. under contract No. DE-AC05-84OR2140 for the U.S. Department of Energy.

DSW

1-
2-
3- **DETECTING ISOTOPIC RATIO OUTLIERS**

4-
5- **Charles K. Bayne¹ and David H. Smith²**

6-
7- ¹**Computing and Telecommunications Division, Oak Ridge National**
8- **Laboratory*, Oak Ridge, Tennessee 37831**

9- ²**Analytical Chemistry Division, Oak Ridge National Laboratory*, Oak**
10- **Ridge, Tennessee 37831**

11-
12-
13- **ABSTRACT**

14-
15- An alternative method is proposed for improving isotopic ratio
16- estimates. This method mathematically models pulse-count data and
17- uses iterative reweighted Poisson regression to estimate model
18- parameters to calculate the isotopic ratios. This computer-
19- oriented approach provides theoretically better methods than
20- conventional techniques to establish error limits and to identify
21- outliers.

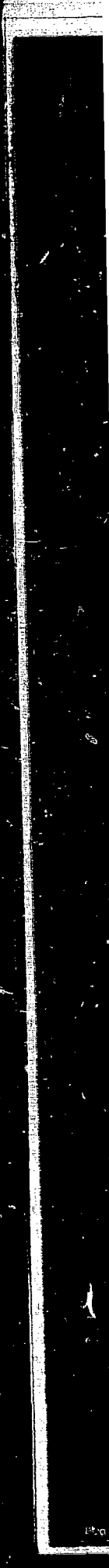
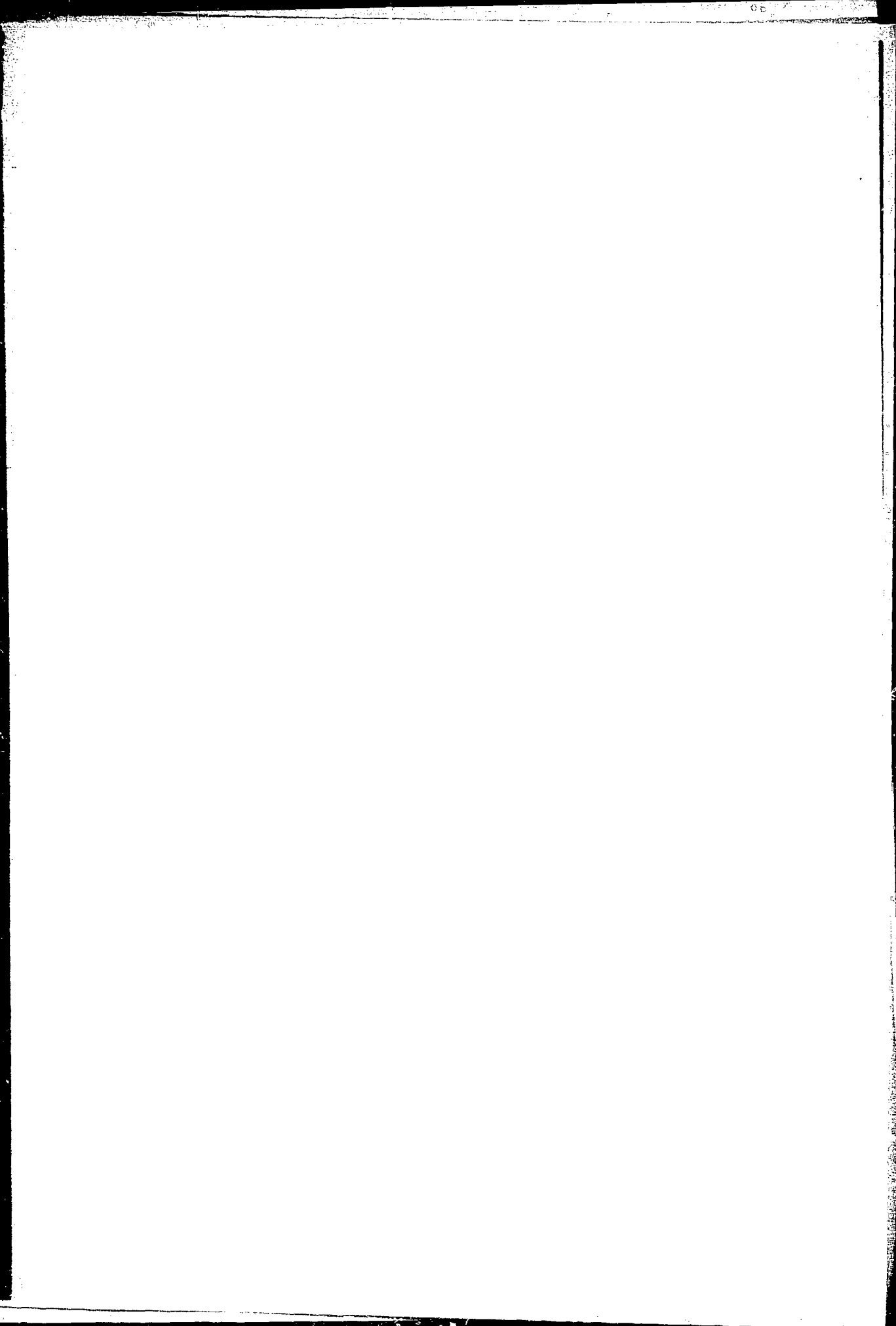
22-
23- **INTRODUCTION**

24- The most sensitive detection technique in mass spectrometry is
25- pulse-counting, where each ion is registered as a count by the
26- detection system. This requires an electron multiplier detector
27- operated at high gain ($>10^6$) and fast amplification electronics to
28- allow reasonable ion counting rates. Most pulse-counting systems
29- have dead times on the order of 10-30 nsec; count rates are
30- generally limited to 5×10^5 /sec or less to reduce count-loss
31- correction due to detector system dead time. Counts are accumulated
32- for each isotope, and calculations of isotopic abundance are made on
33- that basis. Although pulse-counting introduces a few problems not
34- associated with analogue methods, its ability to provide reliable
35- analyses on small samples (1 ng or less for uranium) makes it
36- essential when sample size is limited due to lack of source material
37- or radioactivity.

38- Calculation of isotopic ratios and abundances has previously
39- been executed as if the data were entirely analogous to those
40- generated by current integration techniques. This work presents
41- evidence that, due to their integer nature, such treatment is not
42- proper and introduces a model through which the problem may be
43- addressed.

44-
45-

46- *Operated by Martin Marietta Energy Systems, Inc. under contract No.
47- DE-AC05-84OR2140 for the U.S. Department of Energy.



1- approximated by a normal distribution when the expected values are -48
2- large ($P(\mu)$, $\mu > 30$). However, probability distributions for -49
3- differences and ratios of Poisson counts are not well known. -50
4- -51

5- The problem of deriving the probability distribution for the -52
6- ratio of count data is very difficult. We are not aware of any -53
7- known mathematical expression for this distribution. The range of -54
8- values that the ratio of counts can have is a set of rational -55
9- numbers with many count ratios producing the same rational numbers -56
10- (e.g., $0.5 = 1/2 = 100/200 = 1000/2000 = \dots$). To visualize the -57
11- probability distribution of the ratio of two Poisson variates, -58
12- an example for the ratio of counts from $P(5)$ to counts from $P(25)$ -59
13- was constructed for all values occurring within ± 3 standard -60
14- deviations of the mean for the two distributions. The results are -61
15- given in Fig. 1. From Fig. 1, we note that the distribution of -62
16- count ratios is non-symmetric, skewed to the right, and has more -63
17- than one mode. In addition, values close to the mean, 0.20, can -64
18- have a lower probability of occurring than values further from the -65
19- mean. Values with low probability of occurring that are in the -66
20- middle of the distribution are termed "inliers". This situation -67
21- means that we can encounter "inliers" as well as "outliers" for our -68
22- isotopic ratio estimates. There are no methods presently available -69
23- to detect "inliers". As the means of the Poisson distributions -70
24- increase, we would hope the situation improves. However, for -71
25- expected value ratios of 100:1,000,000, the distribution is still -72
26- non-symmetric, skewed to the right, and the number of modes -73
27- increases. -74

ORNL-DWG 83-19466

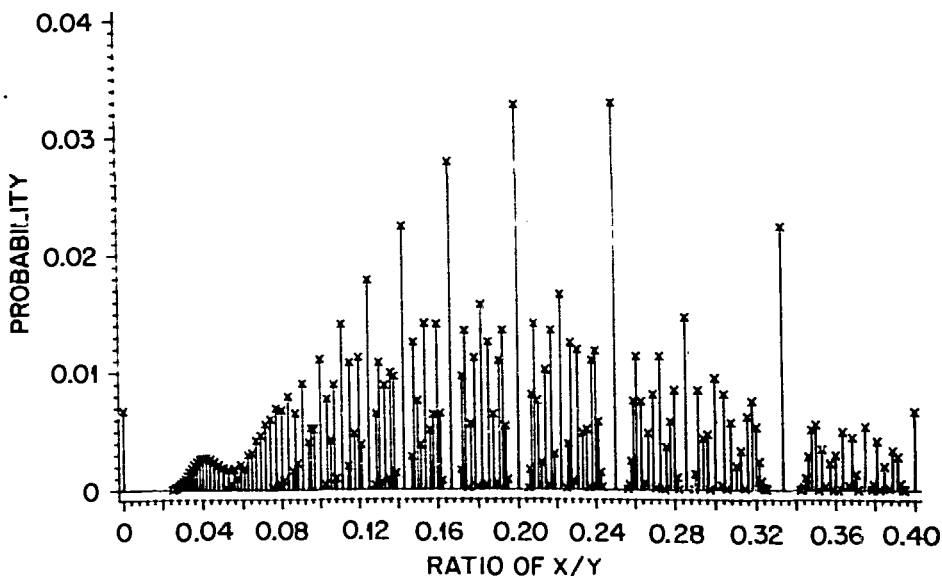


Fig. 1. Probability density function for the ratio of two Poisson variables (X/Y) with expected values of $E(X) = 5$ and $E(Y) = 25$.

1-
2-
3- The greatest handicap with the distribution of ratios of
4- Poisson counts is that the distribution cannot be accurately
5- approximated by normal or other well known probability
6- distributions. This handicap prevents the evaluation of the
7- statistical properties of the isotopic ratio estimates. Chapman [2]
8- showed that the ratio of the expected values from two Poisson
9- probability distributions cannot have an unbiased estimate with a
10- finite variance, although for large counts the bias of the usual
11- estimator is negligible. What is more important is that estimates
12- of precision and percentile points needed to establish confidence
13- intervals of the isotopic estimates cannot be easily obtained. The
14- poor statistical properties of the isotopic ratio estimates
15- necessitate investigating new methods for estimating these ratios.

15- ISOTOPIC RATIO ESTIMATION BY POISSON REGRESSION

16-
17- An alternative method for estimating isotopic ratios, called
18- Poisson regression [3], is based on mathematical models that
19- account for the sources of pulse-count data. These mathematical
20- models contain parameters that represent the isotopic ratios and are
21- estimated by a weighted nonlinear least squares algorithm.

22- The pulse-count data are first adjusted for dead time and sweep
23- factors. The adjusted data are assumed to have a Poisson
24- distribution. That is, if $Y(i,j)$ is the adjusted pulse-count datum
25- in the i -th run and j -th column, its distribution is:

$$26- Y(i,j) \sim P[F(\underline{\theta})] ,$$

27- where $F(\underline{\theta})$ represents the expected value of the pulse-count data as
28- a function of the parameter vector $\underline{\theta}$. This expected value has the
29- following form:

$$31- F(\underline{\theta}) = \text{background} + (\text{bias})^{-1}(\text{run-effect})(\text{isotopic-effect}).$$

32- background = $\exp[B(i)]$, represents background counts.

33- bias = bias factors.

34- run-effect = $\exp[R(i)]$, represents effects that vary from
35- run to run, such as filament-source geometry,
36- sample-filament chemistry, and filament
37- temperature.

38- isotopic-effect = $\exp[I(k)]$, the effect due to the k -th minor
39- isotope where the major isotope has an effect
40- defined to be 1 and the minor isotopes are
41- normalized to the major isotope. In other
42- words, the minor isotopic effect is quanti-
43- tatively a fraction of the abundance of the
44- major isotope.

The mathematical model for the Es example with parameter vector $\theta' = (B(1), \dots, B(10), R(1), \dots, R(10), I)$ and row $i = 1, 2, \dots, 10$ is given by:

$$F[\theta] = \begin{cases} \exp[B(i)], & \text{(background).} \\ \exp[B(i)] + (\text{bias})^{-1} \exp[R(i)], & \text{(major isotope),} \\ \exp[B(i)] + (\text{bias})^{-1} \exp[R(i) + I], & \text{(minor isotope).} \end{cases}$$

This model is nonlinear in the parameters. An iterative nonlinear least squares algorithm is used to solve for the parameter values. Because the variance of a Poisson distribution is equal to its expected value, we weight the data inversely to their variances or

$$W(\theta) = 1/\text{var}(Y(i,j)) = 1/F(\theta) .$$

The estimation criterion is to minimize the weighted sums of squares with respect to the estimated parameter vector, $\hat{\theta}$.

$$S(\hat{\theta}) = \sum_{i,j} W(\hat{\theta}) [Y(i,j) - F[\hat{\theta}]]^2 .$$

To solve for the minimum of the weighted sums of squares, we first linearize the function using a Taylor series approximation. An initial estimate of the parameter vector is made, and a minimum solution can be found using ordinary weighted least squares. The initial solution is updated and the process repeated. Iteration continues until an insignificant change in the model parameter vector occurs. The solution algorithm is shown schematically in Fig. 2. The algorithm used for the solution is a modified algorithm by Frame [4] that is equivalent to finding the maximum likelihood estimates for the parameter vector.

The estimated parameter vector has several asymptotic properties:

1. $\hat{\theta}$ is asymptotically normal.
2. $\hat{\theta}$ is asymptotically unbiased.
3. $\hat{\theta}$ has a lower variance than any other parameter estimate with properties 1 and 2.

These asymptotic properties give a theoretical basis for estimating isotopic ratios and establishing standard deviations and confidence intervals for these estimates.

Two diagnostic tools are available to identify outliers and judge the suitability of the model. The first tool is the Freeman-Tukey (F-T) [5] residuals that are used to identify individual pulse-count values as suspected outlier values. The F-T residuals are standardized to have an approximate normal distribution with zero mean and standard deviation of one. The second tool is the chi-square statistic that is used to judge the appropriateness of removing a suspected outlier or the fit of the model.

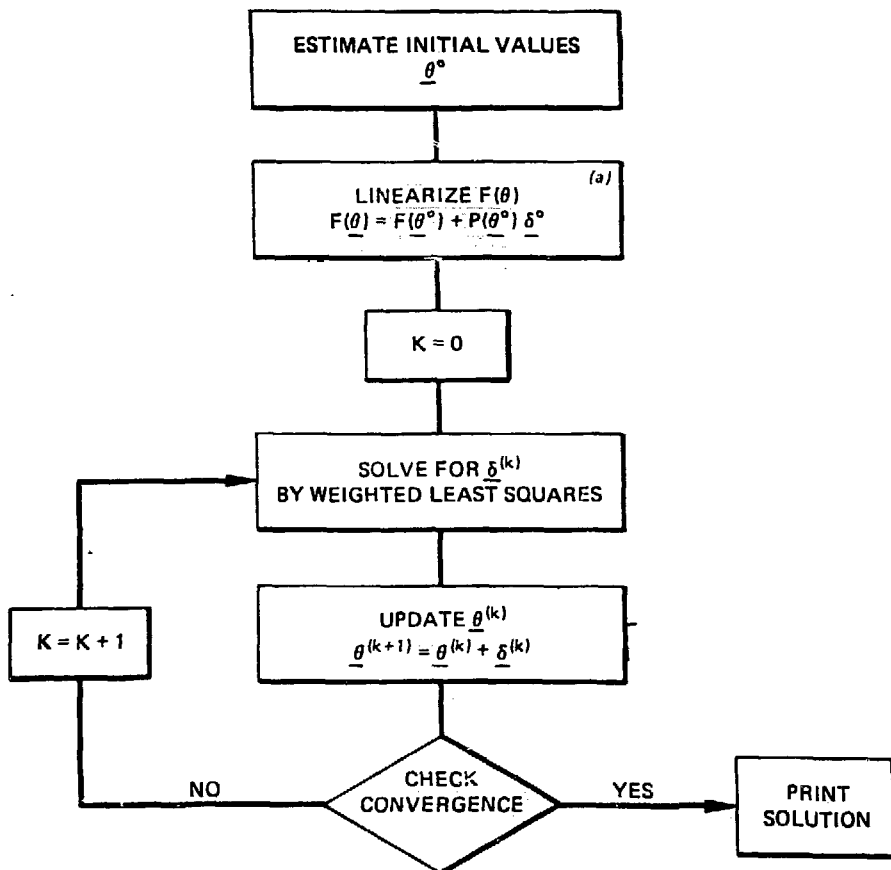
Freeman-Tukey Residual (F-T):

$$F-T = \sqrt{y(1,j) + 1} + \sqrt{y(1,j)} - \sqrt{4F[\hat{\theta}] + 1}$$

Chi-Square Statistic:

$$\chi^2 = \sum_{1,j} \frac{[Y(1,j) - F[\hat{\theta}]]^2}{F[\hat{\theta}]}$$

ORNL-DWG 83-19520A



(a) $P(\theta)$ = MATRIX OF PARTIAL DERIVATIVES OF $F(\theta)$ WITH RESPECT TO θ .
 $\delta^\circ = \theta - \theta^\circ$.

Fig. 2. Schematic diagram of the iterative reweighted least squares algorithm.

An example of these two statistics is given in Table 2 for the einsteinium data. We note that a large value of the F-T residuals occurs in the first row for $^{254}\text{Es}^+$, indicating a potential outlier. By removing this data point and refitting the data, the probability of the chi-square statistic is increased from 0.03 to 0.59. Using the revised data, the 95% confidence interval for the Poisson regression estimate ($^{254}\text{Es}^+ / ^{253}\text{Es}^+ = 0.01965$) is (0.01943, 0.01987).

Table 2. Freeman-Tukey residuals and chi-square statistics for einsteinium data.

Background	$^{253}\text{Es}^+$	$^{254}\text{Es}^+$
1.1	0.8	-3.0 ^a
-0.0	-0.0	0.1
0.1	0.1	-0.4
-0.4	-0.4	1.5
0.3	0.3	-1.0
-0.1	-0.2	0.6
0.0	0.0	-0.1
0.3	-0.2	0.9
-0.5	-0.3	1.2
-0.2	-0.1	0.4

^aSuspected outlier value

Chi-square statistics: (18.3 and 6.5)

1. $\text{Pr}[X^2(9) \geq 18.3] = 0.03$ (with outlier)
2. $\text{Pr}[X^2(8) \geq 6.5] = 0.59$ (without outliers)

The pattern of the F-T residuals can also be used as an indicator for potential physical problems with the sample. For example, variable fractionation of a sample occurs when one or more of the isotopes has a fluctuating volatility rate when compared to other isotopes in the sample. Table 3 gives the usual estimates for a uranium sample of the isotopic ratios $^{234}\text{U}^+ / ^{238}\text{U}^+$, $^{235}\text{U}^+ / ^{238}\text{U}^+$, and $^{236}\text{U}^+ / ^{238}\text{U}^+$. The ten values of each isotopic ratio were checked for outliers using the common outlier tests: Dixon's criterion, studentized deviates, and studentized ranges [6]. Each test indicated that the minimum and maximum isotopic values were not outliers at the 5% significance level. However, an F-T residual plot for the $^{235}\text{U}^+$ and $^{238}\text{U}^+$ counts in Fig. 3 shows large residual for both isotopes for the first three runs and suspected values for $^{238}\text{U}^+$ counts at run numbers four and nine. The residual pattern shows the residuals for the two isotopic counts are about the same magnitude for each run but opposite in sign. This residual pattern has been observed in other samples when variable uranium fractionation was suspected. The value of the chi-square statistic was large ($X^2 = 151.11$, $df = 28$) indicating a significant lack of fit of the model ($P < 0.001$). This result gave additional evidence of physical problems with the sample and indicates it should be reanalyzed.

Table 3. Isotopic ratios for a uranium sample.

Run	$^{234}\text{U}^+ / ^{238}\text{U}^+$	$^{235}\text{U}^+ / ^{238}\text{U}^+$	$^{236}\text{U}^+ / ^{238}\text{U}^+$
1	0.009130	0.9316	0.000983
2	0.009116	0.9342	0.000986
3	0.009144	0.9343	0.001002
4	0.009350	0.9409	0.000964
5	0.009344	0.9395	0.000981
6	0.009356	0.9393	0.000952
7	0.009369	0.9405	0.000977
8	0.009224	0.9394	0.000971
9	0.009257	0.9415	0.001013
10	0.009334	0.9401	0.000964
Average	0.009262	0.9381	0.000979
95% Confidence Interval	± 0.000072	± 0.0025	± 0.000014
Poisson Regression Estimates	0.009291	0.9392	0.000977
95% Confidence Interval	(0.009251, 0.009329)	(0.9387, 0.9396)	(0.000965, 0.000990)

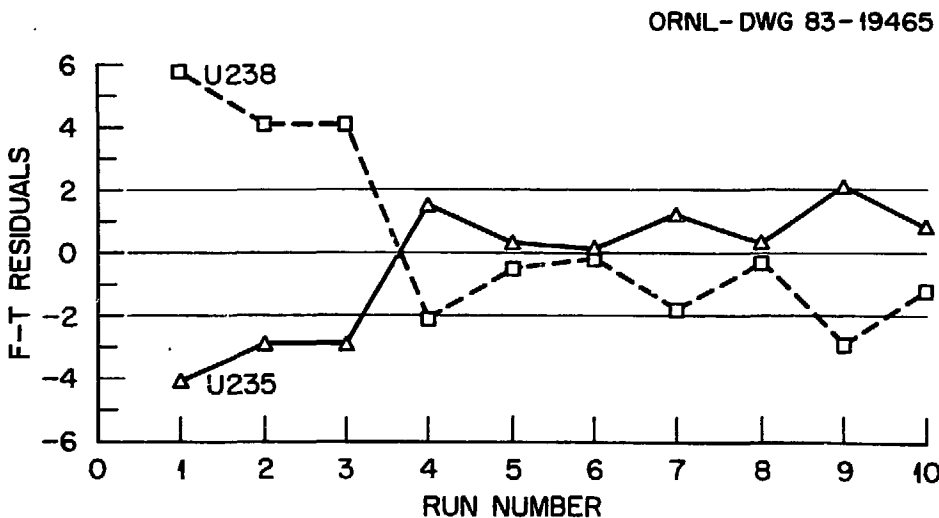


Fig. 3. Freeman-Tukey residuals for $^{238}\text{U}^+$ counts (squares) and $^{235}\text{U}^+$ counts (triangles) from a Poisson regression analysis.

CONCLUSIONS

We have demonstrated that the probability distribution of the ratio of two Poisson variables, such as occurs in pulse-counting mass spectrometry, is ill-defined and certainly not normal. For this reason, it is not appropriate to apply normal statistical treatment to the data. We have developed a model that successfully describes the single-element case. We are currently working on a model to describe more complex cases (isotopic interferences, etc.). The long range goal is to obtain more information concerning the quality of the analysis and thus arrive at better estimates of isotopic ratios.

It is clear that this approach transcends mass spectrometry and is applicable to many analytical techniques where data are collected as individual counts.

REFERENCES

1. E. Serge, Nuclei and Particles, W. A. Benjamin, Inc., New York, 1965.
2. D. G. Chapman, Biometrika, 49 (1952) 45-49.
3. E. L. Frome, M. H. Kutner, and J. J. Beauchamp, J. Am. Stat. Assoc., 68 (1973) 935-940.
4. E. L. Frome, "PREG: A Computer Program for Poisson Regression Analysis", ORAU-178, 1978.
5. Y.M.M. Bishop, S. E. Fienberg, and P. W. Holland, Discrete Multivariate Analysis: Theory and Practice, The MIT Press, Cambridge, Massachusetts, 1975, 130-141.
6. V. Barnett and T. Lewis, Outliers in Statistical Data, John Wiley and Sons, New York, 1978.