# Limitations in the Accuracy of Estimates
# of Blood Group Gene Frequencies

by
*Alexander S. Wiener, M. D.*

The purpose of this communication is to discuss the paper of W. C. Boyd on " The Accuracy of Estimates of MNS Gene Frequencies ". [1] In Boyd's paper, he attempts to demonstrate the advantage of calculations by Fisher's method of maximum likelihood for estimating MNS gene frequencies, as compared with the estimates obtained by certain simple square root formulae used by the present writer. [2] For reasons which will be pointed out, the application of complex maximum likelihood method when estimating gene frequencies represents misdirected energy, since it fails to take into account certain basic limitations of investigations on blood group distribution.

### Significant Figures

The first point raised by Boyd is that of significant figures. As the present writer has pointed out, there is a tendency among workers who use the maximum likelihood method to report gene frequencies with an exorbitant number of significant figures, far exceeding the number justified by the size of the series examined. For example, Fisher [3] has reported gene frequencies with as many as seven significant figures based on tests on fewer than one thousand individuals, while Boyd himself has used five or six significant figures for gene frequencies based on tests on even shorter series. In his new article Boyd admits the validity of this criticism, but at the same time misquotes the present writer as stating that gene frequencies should never be reported with more than two significant figures. The correct quotation is that no more significant figures should be used than are justified by the size of the series. In fact, as Boyd confesses, the present author has reported gene frequencies with three decimal places or significant figures.

By placing excessive emphasis on the use of the complex maximum likelihood calculations, Boyd directs attention away from far more important sources of error in studies on the distribution of the blood groups. Perusal of literature shows many published reports where simple inspection is sufficient to establish the presence of gross errors. The reason for this is that much of the published data is based on tests carried out by workers with only superficial knowledge and understanding of blood group serology, and who are therefore not qualified to carry out accurately

95

the delicate M-N-S and Rh-Hr tests. In fact, there are many reports dealing with the distribution of the basic four A-B-O groups in which blatant errors are manifest, even though the technique of A-B-O testing has become standardized and simplified by the ready availability of potent and specific antisera. For example, the early report of Moss gives the frequencies of the A-B-O groups, based on tests on 80 individuals, as group O, 43.2 percent; group A, 40.0 percent; group B, 7.0 percent; and group AB, 10.0 percent; but gene frequency analysis proves that Moss's frequency for group AB is entirely too high. Whether this was caused by technical or sampling errors is uncertain; at any rate, these results were quoted in textbooks for decades and even at the present time Moss's percentages are occasionally quoted as standard for the U.S.A. A more recent and blatant example is provided by the mass grouping tests carried out on members of the U.S. Armed Forces. Upon re-examination, the A-B-O groups printed on the dog tags proved to be incorrect in as many as 10 percent of the tests. Obviously, technical faults alone could hardly account for such a high percentage of mistakes, which apparently were caused principally by clerical errors.

Another recent example is the report by Gurevitch et al. [4] on the distribution of the Rh-Hr types in certain Jewish populations. In one series of 129 individuals they reported 14 to be of the rare type $Rh_zRh_1$; in a second series of 162 individuals, 9 were reported of the rare type $Rh_zRh_1$; while in a third series of 200 individuals, 6 were reported of the rare type $Rh_zRh_1$. These extraordinary results amost surely represent errors in technique, yet the data were reported without comment by the investigators, nor were they challenged by the editor of the journal in which they appeared. Even experienced investigators may make errors in blood grouping if the reagents they use are too weak or lack specificity. [1] In contrast to sampling errors attributable to the size of the series, technical and clerical errors have a bias. Therefore it is far more important to concentrate one's efforts in order to avoid technical errors which can affect the accuracy of the first and second significant figures, rather than waste valuable time in carrying out complex mathematical computations in order to achieve a dubious increase in " accuracy " affecting the third or fourth significant figure.

## Limitations of the Maximum Likelihood Method

Fisher and Boyd persistently overlook the most serious limitation of the maximum likelihood method, namely, that it entails a knowledge of mathematics which all but a few workers in the field lack, and that the calculations are so complex and laborious that they are costly of time and effort and are subject to error even in expert hands. For example, Fisher [3] devoted an entire article to the complex calculations necessary to derive the Rh-Hr gene frequencies for a single series of individuals tested in England. Obviously such complex calculations would be too impracticable to apply to more than a few of the large number of series given in the book by Mourant. [5] Moreover, in Boyd's article he confesses to certain serious errors in calculations which he himself has made. Thus, the standard deviations for certain gene frequencies $L^s$, $L$, $l^s$, and $l$ were previously given by Boyd as .018, .022, .014, and .020, respectively, but he

96

now confesses too errors in calculations so that these values should have been .024, .027, .021, and .025, respectively, an error of 25 percent. Obviously if an expert mathematician like Boyd can make such gross errors in calculations, the chance of mistakes will be even greater when other less proficient workers attempt to apply the maximum likelihood method. This further justifies the contention of the present author that it is much easier and safer to carry out a few more blood tests, in order to increase the precision of estimated gene frequencies, than to apply the complex maximum likelihood method.

Close examination reveals that the claimed gain in " accuracy " by the use of the maximum likelihood method is too trivial to be worth the cost and effort, since almost identical results can be obtained with the aid of the simple square root formulae. Thus, Boyd confesses that gene frequencies for his Benghali data, as computed by the present author using the simple square root formulae, do not differ from the values obtained by him, using the complex maximum likelihood method. This, he asserts, is accidental, and he cites Fry's Cook Islander data as an " outstanding " example supposedly demonstrating the superiority of the maximum likelihood method. A series of 267 individuals were examined by Fry, so that maximum difference in incidence that could be detected in this series would be of the order of one part in about 250. Boyd reports that the present author's square root method gives a value for one of the gene frequencies of 57.87 percent, while the maximum likelihood method gives the value of 58.04 percent, a difference of about one part in about 500, such as might be obtained by testing one-half of one person in the series of 267. He offers no evidence to show which of the two estimates is closer to the " exact " value, but asserts that the " grater accuracy " of the maximum likelihood estimate is equivalent to testing 17 more Cook Islanders at an increased cost of $ 69.00. Even accepting Boyd's estimate of the cost, the writer wonders how many people would be willing to spend $ 69.00 for an " increased accuracy " in an estimate equivalent to half of one person? Besides, the expert mathematician who makes the maximum likelihood calculations surely is paid a salary! Moreover, as has already been shown, the supposed increase in accuracy of the third significant figure is of dubious value, considering that there are so many other more potent factors affecting the accuracy of the first and second significant figures.

## Comment

During the past few years, a large number of articles have appeared dealing with the distribution of the A-B-O groups and Rh factor in a variety of diseases, such as diabetes, hypertension, gastric and duodenal ulcer, toxemia of pregnancy, carcinoma of the lung and stomach, and sarcoidosis. This repeats a line of investigation which was in vogue more than 30 years ago, and was long ago shown to be both fruitless and pointless. In the more recent studies, small but " statistically significant " differences have been reported in the blood group distribution in various diseases as compared with healthy individuals in so-called " control series ". This has caused these investigators to conclude that individuals of certain blood groups and Rh types

are more susceptible to peptic ulcer, or diabetes or sarcoidosis, etc. They have even gone so far as to speculate how blood group polysaccharides of a particular chemical structure may cause an individual to be more susceptible to certain diseases. These studies, like the use of the maximum likelihood method for calculating blood group gene frequencies, place emphasis on superficialities and disregard more fundamental causes of variations in blood group distribution.

To justify comparing the blood group distribution in a series of individuals with a specific disease and a so-called control series of " normals ", one must make the *tacit assumption that the population from which the individuals were derived is* stable and is homogeneous. However, there is probably no living population today, anywhere in the world, that is strictly homogeneous in a sense that there is random mating among members of the population. Social, economic, political, religious, and other barriers lead to the formation of isolates within the population, so that random mating is the exception rather than the rule. In addition series from different age groups are not comparable because of the effects of migration and differential fertility. In large cities such as New York and London, isolation and inbreeding occur among members of subgroups of different nationalities or religious beliefs, and this is especially true for populations from different cities of the same country. For example, in the U.S.S.R., the percentage of group B has been found to range from 31.5 percent for Asiatic Russians to 46.2 percent in Odessa. If one were to apply the superficial reasoning used by investigators who have been studying blood group distribution in disease, one would conclude that individuals of group B have a special predilection to reside in Odessa in preference to other parts of Russia. Perhaps some percularity of their blood group polysaccharides causes group B people to prefer cities beginning with the letter " O " ?!

Another example of lack of insight in modern blood group literature is provided by the articles in which the incorrect C-D-E notations for the Rh-Hr types are used. The confusion which exists in the literature regarding the serology, genetics and nomenclature of the Rh-Hr types is caused by the fact that most of the articles published have been written by individuals who have only a superficial understanding of the subject. [6, 7, 8] Adopting the simple but incorrect tacit assumption of a one-to-one correspondence between agglutinogen and antibody, such workers merely read the labels of purchased antisera which they use, and then " type " bloods as C+,D—, E+, c—, etc. Certainly, there is much more to Rh-Hr serology and genetics than acquaintance with the English alphabet or the ability to read labels on bottles of antisera. Yet, such individuals write articles on blood grouping, and pose as experts in medicolegal cases of disputed paternity. Recently the present author had the task of correcting errors in M-N and Rh-Hr typing committed by two such self-designated experts, which almost caused serious miscarriages of justice. (Nor does expert knowledge of Chi Square tests make one an expert in blood group serology). Significantly, *the same workers who place emphasis on the supposed value of maximum likelihood* calculations have also endorsed the recent investigations on blood group distribution in disease, and the use of the superficially attractive but fallacious C-D-E notations.

## Summary and Conclusions

1. For reporting results of calculations of blood group frequencies no more significant figures should be used than are justified by the size of the series tested.

2. Technical and clerical mistakes have been responsible for gross discrepancies in reports on blood group distribution published in the literature, which qualified investigators can recognize by mere inspection of the data.

3. Maximum likelihood calculations are complex and laborious and require a mastery of mathematics which few workers possess. Even in the hands of experts gross errors have resulted when carrying out the laborious computations, so that such calculations are not only expensive but impracticable and dangerous.

4. The application of maximum likelihood calculations by Boyd to two series of individuals tested for the M-N-S types has yielded estimates of gene frequencies differing only in the third decimal place from the values obtained by the author with the aid of simple square root formulae. When one considers, even disregarding the possibility of technical and clerical errors, that the sizes of these two series justified the use of only two significant figures, it is evident that the maximum likelihood method contributed only a dubious increase in ,, accuracy ''. Therefore, such calculations consume misdirected energy and place emphasis on superficialities while disregarding more important factors which can affect the accuracy of the data.

5. Many investigators when studying the distribution of the blood groups make the tacit but incorrect assumption that the population is undergoing random mating and is therefore homogeneous. This appears to be one basis for the recent series of fallacious reports on the supposed relationship between the blood groups and various diseases.

6. Confusion in the blood grouping field is due to the fact that most of the articles and books which are published have been written by individuals who have only superficial knowledge and understanding of the subject. The widespread use of the incorrect C-D-E notations reflects the lack of mastery of Rh-Hr serology shared by these workers.

7. It is significant that the same workers who are enthusiastic for maximum likelihood calculations also endorse the pointless studies on blood group distribution in disease, and the use of the superficially attractive but fallacious C-D-E notations.

# References

1. BOYD, W. C.; The « Accuracy » of Estimates of MNS Gene Frequencies. Acta Med. Genet. et Gemellol., 1956.
2. WIENER, A. S.: Serology, Genetics and Nomenclature of the M-N-S Types. Acta Med. Genet. et Gemellol. *3*, 314, 1954.
3. FISHER, R. A.: The Fitting of Gene Frequencies to Data on Rhesus Reactions. Ann. Eugen. (London) *13*, 150, 1954.
4. GUREVITCH, J., and MARGOLIS, E.: Blood Groups in Jews from Iraq. Ann. Human Genet. (London) *19*, 257, 1954-55.
5. MOURANT, A. E.: The Distribution of the Human Blood Groups, 439 pp., C. C. Thomas, Springfield, Ill., 1954.
6. WIENER, A. S., and WEXLER, I. B.: The Interpretation of Blood Group Reactions, with Special Reference to the Serology and Genetics of the Rh-Hr types. Novant'Anni delle Leggi Mendeliane, dell'Istituto « Gregorio Mendel », edited by Prof. L. Gedda, pp. 147-162, Roma, 1956.
7. — The Rh-Hr Types. A Budget of Paradoxes. J. Forensic Med., 2, 224, 1955.
8. — OWEN, R. D., STORMONT, C., and WEXLER, I. B.: Medicolegal Applications of Blood Grouping Tests. J. Amer. Med. Assoc., *161*, 233, 1956.

## RIASSUNTO

1) Per riferire i risultati di calcoli sulle frequenze dei gruppi sanguigni non dovrebbe essere usato un numero di cifre superiore a quello che sia giustificato dall'entità del materiale esaminato.

2) Errori tecnici e di trascrizione sono stati la causa di grossolane discrepanze nelle relazioni sulla distribuzione dei gruppi sanguigni pubblicate nella letteratura, errori che studiosi qualificati possono identificare in base ad un semplice esame dei dati.

3) I calcoli di massima probabilità sono complessi e laboriosi e richiedono una padronanza della matematica che pochi ricercatori possiedono. Anche in mano ad esperti sono risultati degli errori grossolani nell'eseguire questi laboriosi computi, al punto che tali calcoli divengono non soltanto costosi ma anche poco pratici e pericolosi.

## RÉSUMÉ

1) Lorsque l'on rapporte le résultat des calculs faits sur la fréquence des groupes sanguins, on ne doit pas présenter de chiffres plus significatifs que ceux justifiés par le nombre de sujets de la série étudiée.

2) Un investigateur compétent peut reconnaître d'un simple coup d'œil les erreurs techniques et cléricales responsables des variations importantes notées dans la littérature touchant la distribution des groupes sanguins.

3) Les calculs de probabilité maxima sont délicats et difficiles et requièrent une maîtrise des mathématiques que peu de chercheurs possèdent. Même entre mains expertes, des erreurs grossières se glissent au cours de ces calculs laborieux de sorte que de telles méthodes sont non seulement dispensieuses mais inutilisables et dangereuses.

4) L'application par Boyd des calculs de probabilité maxima dans deux séries de cas étudiés au point de vue M.N.S. a conduit à des résultats ne différant qu'à la troisième décimale de ceux obtenus par l'auteur alors que celui-ci n'employait que de simples règles à racine carrée. Si l'on admet, faisant abstraction de la possibilité d'erreurs techniques et cléricales, que l'envergure des deux séries étudiées justifie l'emploi de seulement deux chiffres significatifs, il est évident que la méthode de calcul de probabilité maxima ne fournit qu'un apport de précision négligeable. De tels calculs ne font que consommer une énergie mal dirigée, donner de l'emphase à des problèmes superficiels et négliger par ce fait même des facteurs plus importants qui peuvent influer sur la précision des statistiques.

5) Dans l'étude de la distribution des groupes sanguins, plusieurs investigateurs font l'assomption tacite mais fausse que

4) L'applicazione dei calcoli di massima probabilità, da parte di Boyd, a due serie di individui esaminati per i tipi M-N-S è risultata in valutazione della frequenza dei geni che differiscono soltanto nella terza cifra decimale da quelle ottenute dall'autore con l'aiuto di semplici formule di radici quadrate. Quando si consideri, anche senza tener conto della possibilità di errori tecnici ed umani, che l'entità di queste due serie giustificava l'uso di due sole cifre significative, è evidente che il metodo della massima probabilità contribuiva soltanto ad un dubbio aumento nella « accuratezza ». Perciò tali calcoli consumano energia male sfruttata e mettono in risalto elementi superficiali mentre non tengono conto dei più importanti fattori che possono influenzare l'accuratezza dei dati.

5) Molti ricercatori nello studiare la distribuzione dei gruppi sanguigni suppongono, tacitamente ma erroneamente, che la popolazione proceda ad accoppiamenti indiscriminati e sia perciò omogenea. Sembra che questa sia l'origine della recente serie di erronee comunicazioni sulla supposta relazione fra i gruppi sanguigni e diverse malattie.

6) La confusione nel campo dei gruppi sanguigni è dovuta al fatto che la maggior parte degli articoli e dei libri pubblicati sono stati scritti da persone che hanno soltanto una superficiale conoscenza e comprensione dell'argomento. L'uso assai diffuso della inesatta nomenclatura C-D-E riflette la mancanza di padronanza della serologia Rh-Hr che è comune a queste persone.

7) Vale la pena di notare che gli stessi studiosi che sono entusiasti dei calcoli di massima probabilità appoggiano anche gli inutili studi sulla distribuzione dei gruppi sanguigni nelle malattie e l'uso della nomenclatura C-D-E, superficialmente attraente ma inesatta.

la population se multiplie au hasard et est donc homogène. Ceci semble être une des bases de la série récente d'articles fallacieux sur les supposés rapports entre les groupes sanguins et diverses maladies.

6) La confusion dans le domaine des groupes sanguins est due au fait que la plupart des articles et des volumes publiés ont été écrits par des personnes qui n'ont qu'une connaissance et une compréhension superficielles de ce sujet. Le vaste usage par ces auteurs de l'inadéquate nomenclature CDE reflète une lacune dans leur maîtrise de la sérologie Rh-Hr.

7) Il est remarquable que ces auteurs mêmes, enthousiastes de la méthode de calculs de probabilité maxima, encouragent l'étude inutile de la distribution des groupes sanguins dans les maladies et l'usage de la nomenclature C-D-E attrayante en surface mais fallacieuse.

101

## ZUSAMMENFASSUNG

1. Bei der Berichterstattung der Kalkulierergebnisse mit Bezug auf die Häufigkeit von Blutgruppen sollten keine bedeutsamere Zahlen benützt werden als durch die Grösse der untersuchten Reihenfolgen gerechtfertigt sind.

2. Technische Irrtümer und Schreibfehler müssen für krasse Widersprüche in den in der Literatur feröffentlichten Berichten über die Blutgruppenverteilung · verantwortlich gemacht werden, und derartige Irrtümer sind für den fachkundigen Prüfer durch eine blosse Ueberprüfung leicht erkennbar.

3. Grösstwahrscheinlichkeitrechnungen sind kompliziert und mühselig, und erfordern eine Mathematikkenntnis die nur wenige Facharbeiter besitzen. Grobe Irrtümer ergaben sich sogar in den Arbeiten von Spezialisten bei der Durchführung der mühseligen Berechnungen, so dass solche Kalkulierungen nicht nur teuer sind, sondern auch unpraktisch und gefährlich.

4. Die Anwendung der Grösstwahrscheinlichkeitrechnu n g e n durch Boyd bei zwei Gruppen von Personen die auf die M-N-S Typen geprüft wurden, ergaben Berechnungen von Geneshäufigkeit, die nur an dritter Dezimalstelle von den Zahlen abwichen welche der Autor mit Hilfe einer einfachen Quadratwurzelrechnung erstellte. Wenn man sogar von möglichen technischen irrtümern und Schreibfehlern absieht und in Erwägung zieht, dass die Grösse dieser zwei Reihefolgen die Anwendung von nur zwei bedeutsamen Ziffern gerechtfertigen, ist es klar, dass die Wahrscheinlichkeitrechnung nur einen fragwürdigen Beitrag an « Genauigkeit » darstellt. Derartige Kalkulationen verbrauchen desshalb falsch angewandte Energie und betonen Oberflächlichkeiten, während die wichtigeren Faktoren welche die Genauigkeit der Angaben beeinflussen können unberücksichtigt bleiben.

5. Viele Prüfer gehen beim Studium der Blutgruppenverteilung stillschweigend jedoch unrichtig von der Voraussetzung aus, dass die Bevölkerung sich aufs Geratewohl versippt und desshalb homogen ist. Auf dieser Annahme scheint teilweise eine jüngst aufgetretene Reihe von trügerischen Berichten aufgebaut zu sein die sich mit der mutmasslichen Beziehung zwischen Blutgruppen und verschiedenen Krankheiten befassen.

6. Die Verwirrung auf dem Gebiete der Blutgruppierung ist darauf zurück zu führen, dass der Grossteil der veröffentlichten Artikel und Bücher von Personen stammt die nur oberflächlich mit der Materie vertraut sind. Die weitverbreitete Anwendung der unrichtigen C-D-E Bezeichnungen bezeugt einen Mangel an Rh-Hr Serumkundigkeit welche diesen Arbeitern gemein ist.

7. Es ist bezeichnend, dass dieselben Facharbeiter die mit der Grösstwahrscheinlichkeitrechnung begeistert sind auch die sinnlosen Studien der Blutgruppenverteilung bei Krankheiten befürworten, zusammen mit der oberflächlich anziehenden, jedoch trügerischen C-D-E Bezeichnung.