20th IMEKO TC4 International Symposium and
18th International Workshop on ADC Modelling and Testing
Research on Electric and Electronic Measurement for the Economic Upturn
Benevento, Italy, September 15-17, 2014

# Forensic Metrology: Uncertainty of Measurements in Forensic Analysis

Giuseppe Schirripa Spagnolo[1], Donato Papalillo[1], Lorenzo Cozzella[1], Fabio Leccese[2].

[1.] *Università degli Studi "Roma Tre" – Dipartimento di Matematica e Fisica,*
*Via della Vasca Navale, 84, 00146 Roma, Italy, giuseppe.schirripaspagnolo@uniroma3.it,*
*donato.papalillo@uniroma3.it, lorenzo.cozzella@uniroma3.it.*
[2] *Department of Science, Università degli Studi "Roma Tre", viale Marconi n.84, 00146 Rome,*
*Italy, fabio.leccese@uniroma3.it*

*Abstract* – **In many cases, forensic scientists rely on measurements as a basis for their opinions. In the past, forensic scientists, testifying about such measurements, have often presented the court with a single point value. The problem is that there is an unavoidable and an inherent element of uncertainty in every measurement. The metrology has developed several methods of quantifying a measurement's margin of error or uncertainty. By using these methods, the testifying expert can put the trier in a much better position to determine the appropriate evidentiary weight of the measurement.**
**The purpose of this paper is to describe the problem of uncertainty in forensic science and in particular uncertainty in measurements by forensic scientists.**
**In this paper, the problem of uncertainty is described in a relatively non-technical way. In other words, the discussion will focus on concepts rather than on analytical analysis.**

## I. INTRODUCTION

Forensic metrology is the application of measurements and measurement standards to the solution and prevention of crime [1].

Forensic scientists play a pivotal role in the criminal justice system, providing crucial information about the evidence to the trier of fact. Because the work they do both at the crime scene and in the laboratory often must be used in court, it is especially important that

The goal of forensic science is to provide to justice system an unbiased, independent scientific analysis with expert testimony. When done well, forensic science helps investigators into gather evidence from crime scenes and catch suspects, and it helps the courts correctly determine guilt or innocence.

Popular television programs, like CSI, give the idea that scientific techniques can link a mark (e.g., a partial fingerprint or tireprint) to a unique unknown source. This leads, in the collective public imagination, to the conclusion that supposed scientists have the possibility to

certainly identify the gun that fired the murderous bullet.

It is natural to ask: Forensic science really makes such pinpoint determinations? Can forensic scientists be sure that a particular hammer, excluding all other hammers in the world, produced the imprints observed on a victim's body? The concept of individualization which lies at the core of numerous forensic science subfields exists only in a metaphysical or rhetorical sense. There is no scientific basis for such individualization claims.

Unfortunately, these misconceptions are rooted in many tries and lawyers.

Forensic scientists are not able to link a fingerprint, a hair, a handwriting sample, a tiremark, a toolmark, or any other evidentiary forensic item to its unique source, even though they assert such ability every day in courts. The issue is not the sincerity of the beliefs of workaday forensic scientists. The issue is whether any scientific evidence, that can support those beliefs, exists. There is no theoretical basis or data to the core thesis that every distinct object leaves its own unique set of markers that can be identified by a skilled forensic scientist. Their claims exaggerate the state of their science.

With the importance of forensic science to truth and justice the science employed, relied upon by judges and juries must be valid. Furthermore, the forensic scientists must use reliable techniques, use accepted measurement protocols and provide the limits and uncertainties of the tests.

*"A measurement result is complete only when accompanied by a quantitative statement of its uncertainty. The uncertainty is required in order to decide if the result is adequate for its intended purpose and to ascertain if it is consistent with other similar results."*[2].

Knowledge of the uncertainty associated with measurement results is essential to the interpretation of the results themselves. Without assessments of uncertainty, it is impossible to decide whether observed differences between results reflect more than experimental variability, or whether test items are comply with specifications, or whether laws based on limits have been broken. Without

information on uncertainty, there is a risk of misinterpretation of results. Incorrect decisions taken on such a basis may result in unnecessary expenditure in industry, incorrect prosecution in law, or adverse health or social consequences.

It does not matter how well forensic scientists abide by testing protocols, or how reliable the techniques are, if the underlying science does not actually reveal what the expert says it does. Method validation studies and new research must be on-going even in the area of traditional forensic disciplines.

The uncertainty associated with forensic scientific investigation it is an emerging Branch in Metrology [3].

The uncertainty measurement should be included for both quantitative (how much is there) and qualitative (is it there) analysis [4].

## II.  THE FORENSIC CONTEXT

An expert fulfils a unique function in Court by assisting the Trier of Fact to understand technical matters not commonly known to a lay person. Exclusion rules can be used for determining whether the evidence is admissible at trial or not. In essence, to be admissible, the expert evidence must be relevant to the issues at trial, helpful in assisting the decision, beyond common knowledge and within the expert's sphere of expertise. To be admissible, expert opinion must be based upon proven facts, results of expert investigation and analysis. Virtually anyone can be an expert by suitably qualifying themselves. The extent of the expertise so established ultimately determines the weight given to the expert's evidence by the Trier of Fact. The expert is responsible to the Court in that the evidence given must be complete, impartial and unbiased.

Some key guidelines were listed below, to assist in the interpretation of this Exclusion Rule when applied to scientific knowledge. The guidelines are:
a) whether the theory or technique could be or has been tested.
b) whether the technique has been published or subject to peer review.
c) whether the actual or potential error rates have been considered.
d) whether the technique is widely accepted within the relevant scientific community.

Clearly, these guidelines give much more specific definition to the criteria that should be applied by the courts when considering expert evidence, and allow much more opportunity for new methods that could be utilized.

Improperly, in the Forensic metrology there are two different points of view, namely a "point" paradigm and a "set" paradigm as discussed below [5].

The point paradigm is characterized by the notion that each measurement results in a single, "point-like" value which could in principle be the true value. As a consequence each measurement is independent of the others and the individual measurements are not combined

in any way. In its most extreme form, this way of thinking manifests itself in the belief that only one single measurement is required to establish the true value [6].

The set paradigm is characterized by the notion that each measurement is only an approximation to the true value and the deviation from the true value is random. As a consequence, a number of measurements are required to form a distribution that clusters around some particular value. The best information regarding the true value is obtained by combining the measurements using theoretical constructs in order to describe the data collectively. The operational tools that are available for this purpose include the formal mathematical procedures that can be used to characterize the set as a whole, such as calculating the mean and the standard deviation.

Although it may seem obvious, the correct procedure is the set paradigm, *"forensic science should be scientific"*. Unfortunately, forensic scientists often reject error rate estimates in favor of arguments that describe theirs science error-free. For example, an FBI document section chief asserted that all certified document examiners in the United States would agree with his conclusions in every case [7].

*Table 1: **Point paradigm vs. Set paradigm***

| Point paradigm | Set paradigm |
|---|---|
| The measurement process allows you to determine the true value of the measurand. | The measurement process provides incomplete information about the measurand. |
| "Errors" associated with the measurement process may be reduced to zero. | All measurements are subject to uncertainties that cannot be reduced to zero. |
| Each single reading is potentially the true value of the measurand. | All available data are used to construct distributions from which the best approximation of the measurand and an interval of uncertainty are derived. |

## III.  UNCERTAINTY IN QUANTITATIVE MEASUREMENTS

Testing laboratories shall have and apply procedures for estimating uncertainty of measurement [8,9]. In certain cases the nature of the test method may preclude rigorous, metrologically and statistically valid calculation of uncertainty of measurement. In these cases the laboratory shall at least attempt to identify all the components of uncertainty and make a reasonable estimation, and shall ensure that the form of reporting of the result does not give a wrong impression of the uncertainty.

Reasonable estimation shall be based on knowledge of the performance of the method and on the measurement scope and shall make use of, for example, previous

experience and validation data.

When you make a measurement, there are some questions or uncertainty about how the measured value refers to the true value of the measurand.

When a measurement is made, there are some questions or uncertainty as to how the measured value relates to the true value of the measurand. The measuring system, measurement procedure, operator skill, environment, or other factors may introduce random and/or systematic errors into the measuring process. In general, any series of replicated measurements (or sample) will result in a dispersion or distribution of measured values with the true value of the measurand lying somewhere within the limits of the measurement results.

For a normally or uniformly distributed sample of replicate measurements, the best estimate of the true value of the measurand is given by the sample mean or average. The best estimate of the dispersion of the measurements is given by the standard deviation of the sample.

A generally accepted approach to estimating the measurement uncertainty involves calculating an expanded interval with an associated coverage probability. The expanded interval is determined by:

✓ calculating or estimating the variances of the sources of uncertainty;
✓ calculating the combined uncertainty using the root sum square method;
✓ selecting a coverage factor and multiplying the combined uncertainty to obtain an expanded uncertainty;
✓ adding and subtracting the expanded uncertainty to the mean (or measured value) of the measurand to yield, respectively, the upper and lower limits of the expanded interval;
✓ determining the level of confidence of the expanded interval based upon the coverage factor.

Consideration of possible sources/components of uncertainty included: the methods and materials used, environmental conditions, properties and condition of the test item (i.e. the target shotgun), and the operators (analysts). Methods and materials were standardized by means of Instructions to be followed during the exercise. Environmental conditions were taken into account by making repeated measurements at various times over a period of days. The same test item was used by all participants (analysts). The physical properties and condition of the test item was stable and not a source of uncertainty. The operators (analysts) influence was taken into account by making replicate measurements of the measurand.

## IV.  UNCERTAINTY IN QUALITATIVE TESTING

Qualitative testing generally relates to categorical statements, such as "present/absent", "pass/fail", perhaps membership of a class of compounds. Such classification statements are not usually associated with a range of expression; one does not, in general reporting, generally speak of an artefact or material being a 90% pass, or 99% present. Partial class membership is used extensively in "fuzzy logic" systems, but the relevant terminology and treatment is very rare in ordinary testing activities [10].

The identification of the culprit of a crime has always represented a challenge, not only for the police, the jury and the judge, but also for the general public [3].

In order to do this the materials of interest must first be detected. The ability of a technical method to detect a target material depends upon the amount of the material which is present in the system as well as upon the performance characteristics of the method. Thus, if the aim of an analysis is to determine whether or not a particular substance is present, it will be necessary to specify a minimum concentration which must be capable of detection.

Just as it is possible to make an erroneous identification of a person under poor observation conditions so too is it possible to make an erroneous identification of a material submitted for qualitative analysis. It is hence desirable to provide users of qualitative analysis results with some indication of the reliability of identification.

The degree of confidence in the correctness of identification can be expressed in a number of ways. For a given test method, the basic properties that need to be measured are the numbers of true positive and negative results and the numbers of false positive and negative results obtained on a range of samples [10]. From these numbers the fundamental measures of reliability *viz.* the false positive and false negative rates can be calculated. Several other measures can also be derived from these numbers (see Table 1) [11]. The false positive and negative rates can be combined into a single figure expressed by the Bayesian likelihood ratio. If the analyst is able to quantify his initial degree of belief in the outcome of a test applied to a particular sample - before the test is applied - then a further reliability measure in the form of a Bayesian posterior probability can be calculated. One other important method parameter which needs to be determined is the limit of detection; knowing this enables the analyst to select a method capable of satisfying the customer's requirement relating to minimum detectable amount.

In Table 1 the terms sensitivity is the fraction of true positive results obtained when a test is applied to positive samples (it is the probability that a positive sample is identified as such); specificity is the fraction of true negative results obtained when a test is applied to negative samples (it is the probability that a negative sample is identified as such).

Because false response rates are, in general, low for effective methods, it often takes an extremely large number of experiments to obtain even indicative values. Further, observed false response rates are influenced very considerably by the characteristics of the test population. For example, false response rates are much higher when

the typical level of a material falls close to the response threshold of a simple spot test. Thus, it is unrealistic to expect great reliability in false response rates obtained within a laboratory; it is often difficult to obtain false response rate figures accurate to within an order of magnitude.

*Table 1. False positive and false negative rates.*

| Reliability Measure | Expression |
|---|---|
| False positive rate | $\dfrac{FP}{TN+FP}$ |
| False negative rate | $\dfrac{FN}{TP+FN}$ |
| Sensitivity | $\dfrac{TP}{TP+FN}$ |
| Specificity | $\dfrac{TN}{TN+FP}$ |
| Efficiency | $\dfrac{TP+TN}{TP+TN+FP+FN}$ |
| Youden Index | Sensitivity + Specificity - 100 |
| Likelihood Ratio | $\dfrac{1-\text{False negative rate}}{\text{False positive rate}}$ |
| Bayes posterior probability | Bayes rule |

TP = number of true positives;  FP = number of false positives;
TN = number of true negatives; FN = number of false negatives.

In this paper, the following definitions are used:

*True positive*: Results obtained using the confirmatory technique and another analytical technique are both positive.

*True negative*: Results obtained using the confirmatory technique and another analytical technique are both negative.

*False positive*: Result obtained using the confirmatory technique is negative but that obtained using another analytical technique is positive.

*False negative*: Result obtained using the confirmatory technique is positive but that obtained using another analytical technique is negative.

Quantitative expression and reporting of qualitative testing uncertainties is accordingly unlikely to give indicative, but not very accurate information.

Probably the most important alternative to simple statements of false response rates is the use of values derived from Bayes' theorem.

Examples include likelihood ratio (an indication of the

additional information provided by a test result) [12,13] and posterior probability, an indication of the probability of an object fitting a given category given a test result. Bayesian estimates are particularly widely used in evaluating forensic evidence, for example DNA matching or blood group matching. Bayesian estimates can be calculated by appropriate combination of false positive and false negative rates [14,15].

In the likelihood-ratio framework the task of the forensic scientist is to determine the probability of obtaining the observed properties of the sample of known origin and the sample of questioned origin under the hypothesis that the two samples have the same origin versus under the hypothesis that they have different origins, see Eq. (1):

$$LR = \frac{p\left(E|H_{so}\right)}{p\left(E|H_{do}\right)} \tag{1}$$

where *LR* is the likelihood ratio, *E* is the evidence, i.e., the properties of the sample of known origin and the properties of the sample of questioned origin, $H_{so}$ is the same-origin hypothesis and $H_{do}$ is the different-origin hypothesis.

$p(E|H_{so}) =$  Probability of evidence and known (suspect) under Prosecution hypothesis;

$p(E|H_{do}) =$  Probability of evidence and known (suspect) under Defense hypothesis;

A likelihood ratio greater than one lends support to the same origin hypothesis, e.g., if the likelihood ratio is 100, then, whatever the trier of fact's prior belief, after hearing this they should be 100 times more likely than before to believe that the samples have the same origin rather than different origins. Similarly, a likelihood ratio less than one lends support to the different-origin hypothesis, e.g., if the likelihood ratio is 1/100, then, whatever the trier of fact's prior belief, after hearing this they should be 100 times more likely than before to believe that the samples have different origins rather than the same origin. The deviation of the likelihood ratio from one is a quantification of the strength of the evidence with respect to the competing same-origin and different-origin hypotheses.

Decisions are binary, as shown in Table 2. Errors are counted if the system declares two samples to have the same origin when in fact they have a different origin (false positive), or when it declares two samples to have different origins when in fact they have the same origin (false negative). However, as will be explained below, such a metric of validity is not appropriate if one is working within the likelihood-ratio framework.

*Table 2: Correct classifications and classification errors*

| Truth | System output | |
|---|---|---|
| | Same-origin | Different-origin |
| Same-origin | True positive | False negative |
| Different-origin | False positive | True positive |

Correct-classification rates/classification-error rates are not appropriate for use within the likelihood-ratio framework because they are based on posterior probabilities rather than likelihood ratios, and they are based on a categorical thresholding, error versus not-error, rather than a gradient strength of evidence.

In order to calculate a posterior probability, one would have to combine prior probabilities with the likelihood ratio, as shown in Eq. (2) (odds form of Bayes' Theorem).

$$\frac{p(H_{so}|E)}{p(H_{do}|E)} = \frac{p(E|H_{so})}{p(E|H_{do})} \cdot \frac{p(H_{so})}{p(H_{do})} \qquad (2)$$

$$\underset{\substack{posterior \\ odds}}{} \quad \underset{\substack{likelihood \\ ratio}}{} \quad \underset{\substack{prior \\ odds}}{}$$

The prior odds, the trier of fact's belief about the relative probability of the competing hypotheses before the evidence is presented, and the posterior odds, the trier of fact's belief about the relative probability of the competing hypotheses after the evidence has been presented, are not within the purview of the forensic scientist. In fact it is strength of the likelihood-ratio framework that the forensic scientist focuses only on calculating the strength of evidence on the basis of the samples presented to them for analysis, and does not consider the prior odds, thus reducing the likelihood that their analysis could be consciously or unconsciously influenced by other knowledge about the case. Since posterior probabilities are not calculated as part of the likelihood-ratio framework, a metric of validity cannot be based on posterior probabilities; rather it should be based on likelihood ratios of evidence on the basis of the samples presented to them for analysis, and does not consider the prior odds, thus reducing the likelihood that their analysis could be consciously or unconsciously influenced by other knowledge about the case (such potential sources of human bias [2]). Since posterior probabilities are not calculated as part of the likelihood-ratio framework, a metric of validity cannot be based on posterior probabilities; rather it should be based on likelihood ratios.

The size of a likelihood ratio indicates the strength of its support for one hypothesis over the other. It would be worse to report a likelihood ratio of one million in favor of a contrary-to-fact hypothesis (a likelihood-ratio which supports the same-origin hypothesis when the objects actually have different origins, or a likelihood-ratio which supports the different-origin hypothesis when the objects actually have the same origin) than to report a likelihood ratio of ten in favor of a contrary-to-fact hypothesis, because the former provides greater support for the contrary-to-fact hypothesis and would thus have a greater potential to contribute to a miscarriage of justice. A metric of validity should therefore assign a greater penalty to the former than to the latter. Such a gradient metric contrasts with the binary classification-error rate metric which would assign an equal penalty to both.

## V. OBTAING FALSE POSITIVE/NEGATIVE RATES

There are basically two ways of obtaining false response rates for a given technique. The first involves a review of the literature for the particular technique to see if studies on the false response rates have already been carried out and recorded. For techniques commonly used, this information might be expected to be in the public domain. For in-house methods the information should have been generated during method validation studies. Published false response rates should be used with caution; they will have been obtained using particular equipment, and will refer to particular samples; it is necessary consider whether if the situations are comparable.

If information on the false response rates for a particular technique is not available it will have to be generated by an experimental study.

Two mechanisms can contribute towards the production of false responses. In the first of these, false responses are caused by the sample "effects".

One or more components of the sample can interact with the detection system to produce a false positive response. Similarly, one or more components of a sample, other than the target component, can interact with the detection system to inhibit the production of a genuine positive response thereby leading to a false negative response.

A second mechanism can operate near the cut-off region of a test. Here, the number of false positives depends upon the distribution of values obtained on blanks. A cut-off value is selected - typically at a level of 3 standard deviations of the blank - below which values are regarded as negative and above which they are regarded as positive. Thus (for a 3 standard deviations cut-off) there is an a priori probability of obtaining 1 or 2 false positive results in every 1000 tests on genuinely negative samples.

Raising the cut-off level reduce the probability of obtaining false positive but increases that of obtaining false negatives, and conversely. These ideas are illustrated in Figure 1.
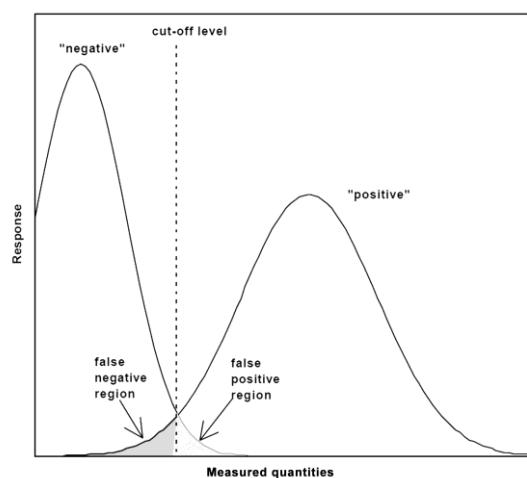


*Figure 1: False response rates from distributions*

Estimation of the false response rates of a method should ideally be designed into the method validation studies. At this stage the technique would of course be known but a study should ensure that an adequate range of samples, likely to be encountered in practice, is covered. A confirmatory detection technique will also need to be selected and a method incorporating it validated. Given that the number of false responses should ideally be low, the problem arises of how many samples to test to be reasonably sure of finding a non-zero number of false responses.

The open problems are the correct determination of the Figure 1 curves and the uncertainty associated with them.

## VI. CONCLUSION

In this paper we have illustrated the problems associated with forensic measures. Many problems still remain open. In any case, forensic scientists must adopt protocols, such as blind examinations in combination with realistic samples that minimize the risks that their success rates will be inflated and their conclusions biased by extraneous evidence and assumptions.

Although obstacles exist both inside and outside forensic science, the time is ripe for the traditional forensic sciences to replace antiquated assumptions of uniqueness and perfection with a more defensible empirical and probabilistic foundation.

## REFERENCES

[1] Ted Vosk, A.F. Emery, "Forensic Metrology: Scientific Measurement and Inference for Lawyers, Judges and Criminalists", CRC Press - Taylor & Francis Group, Boca Raton, FL USA, 2014. ISBN 9781439826195

[2] B.N. Taylor, C.E. Kuyatt, "NIST Technical Note 1297, 1994 Edition, Guidelines for Evaluating and Expressing the Uncertainty", National Institute of Standards and Technology Gaithersburg, MD, USA. available online: http://physics.nist.gov/Pubs/guidelines/TN1297/tn1297s.pdf

[3] A. Ferrero and V. Sotti, "Forensic Metrology: A New Application Field for Measurement Experts Across Techniques and Ethics", IEEE Instrum. Meas. Mag., vol. 9, no. 3, February 2013, pp. 44–51. doi: 10.1109/MIM.2013.6417051.

[4] S.L.R. Ellison, "Uncertainties in qualitative testing and analysis", Accred. Qual. Assur., vol. 5, Issue 8 , August 2000, pp 346-348. doi: 10.1007/s007690000212

[5] A. Buffler, S. Allie, F. Lubben, B. Campbell, "The development of first year physics students' ideas about measurement in terms of point and set paradigms", Int. J. Sci. Educ. vol. 23, Issue 8, November 2001, pp. 1137–56. doi: 10.1080/09500690110039567

[6] M.G. Séré, R. Journeaux, C. Larcher, "Learning the statistical analysis of measurement error", Int. J. Sci. Educ., vol. 15, Issue 4, Jul-Aug 1993, pp. 427-438. doi:10.1080/0950069930150406

[7] M.J. Saks, J.J. Koehler, The coming paradigm shift in forensic identification science, Science vol. 309, no. 5736, August 2005, pp. 892–895. doi:10.1126/science.1111565.

[8] ISO/IEC 17025, "General requirements for the competence of testing and calibration laboratories", Second edition 2005-05-15.

[9] BIPM, IEC, IFCC, ISO, IUPAC, IUPAP, OIML JCGM 100:2008, Evaluation of measurement data - Guide to the expression of uncertainty in measurement (GUM 1995 with minor corrections) (www.bipm.org). (Printed as ISO/IEC Guide 98-3:2008. (www.iso.org))

[10] S.L.R. Ellison, S. Gregory, PerspectiveQuantifying uncertainty in qualitative analysis uncertainty in qualitative analysis, Analyst vol. 123, Issue 5, 1998 , pp. 1155–1161. doi: 10.1039/A707970B

[11] QAWG/03/06, EURACHEM/CITAC Guide: The Expression of Uncertainty in Qualitative Testing, September 2003. Available online: http://www.nmschembio.org.uk/dm_documents/lgcvam2003048_cojjl.pdf

[12] B. Robertson, G.A. Vignaux, "Interpreting Evidence", Wiley, Chichester, UK, 1995. ISBN: 978-0471960263

[13] C.G.G. Aitken, F. Taroni, "Statistics and the Evaluation of Forensic Evidence for Forensic Scientist", 2nd ed., Wiley, Chichester, UK, 2004. ISBN: 9780470011232..

[14] F. Taroni, S. Bozza, A. Biedermann, P. Garbolino, C. Aitken, "Data analysis in forensic science – A Bayesian decision perspective", Wiley, Chichester, UK, 2010. ISBN: 9780470998359.

[15] G. S. Morrison, Measuring the validity and reliability of forensic likelihood-ratio systems, Science and Justice 51, 2011, pp. 91–98. doi: 10.1016/j.scijus.2011.03.002.