

# Bounds on the Mean Classification Error Rate of Multiple Experts

Padhraic Smyth  
Information and Computer Science  
University of California, Irvine  
CA 92717-3425, USA.  
email:smyth@ics.uci.edu

October 29, 1996

## Abstract

A database contains  $N$  items, each item belonging to one and only one of a finite set of classes. The *true* class labels for these items are unknown.  $K$  experts each provide a set of  $N$  classification labels for the  $N$  items in the database. In this paper it is shown that given the experts' labels, one can compute simple bounds on the average classification accuracy of the experts relative to the unknown true labels. No assumptions are made about the labelling patterns of the experts or the nature of the data. The bounds are useful in practical classification problems where absolute ground truth is unknown and experts must subjectively provide labels for feature data. The method is applied to the problem of assessing the collective accuracy of geologists who count volcanoes in images of Venus.

## Keywords:

classification, multiple experts, remote-sensing, planetary geology.

## 1 Introduction and Notation

Consider that a person (an observer) has a database of  $N$  items, each described by a feature vector  $\underline{x}_i$ ,  $1 \leq i \leq N$ . Each item belongs to one of  $m$  classes,  $m \geq 2$ : the classes are mutually exclusive and exhaustive. It is assumed that for each item  $\underline{x}_i$  there exists a true label  $\omega_i$  (a reference label) which is unknown. For example, if the  $\underline{x}_i$  were pixel measurements of an object of unknown class in a remotely-sensed image, the true class label could in principle be obtained by visiting the ground site and ascertaining the class of the object in an unambiguous manner (so called "ground truth")

The observer is assumed to have no information whatsoever about the true class labels of the items. Let  $K$  experts ( $K > 1$ ) each provide a set of  $N$  labels for the  $N$  items, i.e., each expert examines each item  $\underline{x}_i$  in turn and provides a subjective estimate of the true class label for that item. Define  $\bar{e}$  as *the mean classification error rate, averaged across the  $K$  experts, relative to the*

*true labels*, i.e., over all the experts, a certain fraction of items have been mislabelled relative to the truth. By definition

$$\bar{e} = \frac{1}{KN} \sum_{k=1}^K \sum_{i=1}^N e_{ik} \quad (1)$$

where, for the label of labeller  $k$  on the  $i$ th item,  $e_{ik} = 1$  if it is in error and  $e_{ik} = 0$  if it is correct.

The fact that  $\bar{e}$  is defined as the *mean* error rate of  $K$  labellers rather than the error rate of any one labeller is a key point and enables calculation of the bounds. Without knowing ground truth one can not make any statements about the errors of an individual labeller. References to “errors” will be assumed to mean “errors relative to ground truth” throughout the paper.

## 2 Motivation and Background for this Problem

Assessing the collective classification accuracy of a group of experts on a database is an important issue in certain practical classification problems. For example, scientists subjectively label pixels or regions in a remote-sensing image into a set of known ground-cover classes, or medical specialists classify medical records into particular diagnostic classes. In such cases obtaining the true class labels for the data is frequently either physically impossible or prohibitively expensive. For example, in remote-sensing it may be impractical to visit the remote sites to ascertain ground truth. In medical diagnosis it may be too expensive to perform the necessary tests or surgery to determine with absolute certainty what disease the patient actually had. In classification-oriented applications, as online data become more readily available, the proportion of the data for which the true class labels are known is likely to continue to decrease. Quantitative statements about the accuracy of human experts, such as the bounds derived here, are quite valuable in these types of problems. In Section 4 we describe a particular application of the method to counting volcanoes in radar images of Venus. The volcano counting problem originally motivated this work: it is a problem of considerable geologic importance involving multiple expert opinions.

## 3 Bounds on $\bar{e}$

### 3.1 A General Lower Bound on $\bar{e}$

From Equation (1), the average error rate can be written as

$$\bar{e} = \frac{1}{KN} \sum_{i=1}^N e_i \quad (2)$$

where  $e_i = \sum_{k=1}^K e_{ik}$ , is the total number of errors made on item  $i$ ,  $0 \leq e_i \leq K$ .

Consider the  $i$ th item. Let  $n_{ij}$  be the number of times that label  $j$  was provided by the  $K$  labellers for item  $i$ ,  $0 \leq n_{ij} \leq K$ .

Let  $j^*$  indicate the correct label for the item. Thus  $K - n_{ij^*}$  is the number of errors made on

the  $i$ th item. Since  $j^*$  is unknown, one has

$$e_i \geq \min_j \{K - n_{ij}\}, \quad 1 \leq j \leq M. \quad (3)$$

Thus,

$$\bar{e} \geq \frac{1}{KN} \sum_{i=1}^N \min_j \{K - n_{ij}\} = \frac{1}{KN} \sum_{i=1}^N \left( K - \max_j \{n_{ij}\} \right). \quad (4)$$

This bound is a function of the number of disagreements made by the labellers. If there are no disagreements, the bound is 0. In the worst-case scenario the labellers agree on all items but are incorrect in each case, yielding a lower bound of 0 while the true error rate is 1. In general, however, the bound will be non-zero for practical problems, thus providing an indication of the overall error rate of a set of experts. Note that at least one of the  $K$  labellers must have an error-rate *greater than or equal to* the lower bound. Thus, for example, even if the labellers are considered experts in their field, the bound will imply that at least one of them has an error rate greater than some value, relative to ground truth. If this value is large (say greater than 10%) it may indicate the need to re-evaluate the quality of the feature data  $\underline{x}_i$ , or the quality of the expert labelling process, or both.

Equation (4) is the lowest bound one can obtain on the mean error rate without additional information about the problem being available. For example, if  $K = 2$  and one of the labellers is always correct, then the bound is exactly the mean error rate.

### 3.2 A Lower Bound for Binary Classification

With binary classification,  $m = 2$ , we can index the labelling patterns by the number of labels belonging to one of the classes (“detections”),  $0 \leq d \leq K$ . Let  $n_d$  be the number of labelling patterns which have  $d$  detections ( $\sum_{d=0}^K n_d = N$  if all items are labelled). For example,  $n_1$  is the number of items each of which were labelled as a detection by only one of the  $K$  labellers. For binary labels, the bound reduces to:

$$\bar{e} \geq \frac{1}{KN} \sum_{d=0}^K n_d \left( K - \max_j \{n_{ij}\} \right) \quad (5)$$

$$= \frac{1}{KN} \sum_{d=0}^K n_d \min\{K - d, d\} = \frac{1}{KN} \sum_{d=1}^{K-1} n_d \min\{K - d, d\}. \quad (6)$$

With  $K = 2$ , i.e., two labellers,

$$\bar{e} \geq \frac{n_1}{2N} \quad (7)$$

where  $n_1$  is the number of items labelled by the 2 experts where they disagree, i.e., one gets the simple result that the mean error rate is lower bounded by half the fraction of disagreements. If two labellers disagree on all items, their mean error rate must be 0.5 (which also equals the bound in this case).

### 3.3 An Upper Bound on $\bar{e}$

One can also derive a simple upper bound on  $\bar{e}$ :

$$\bar{e} \leq \frac{1}{KN} \sum_{i=1}^N \max_j \{K - n_{ij}\} \quad (8)$$

$$= \frac{1}{KN} \sum_{i=1}^N \left( K - \min_j \{n_{ij}\} \right). \quad (9)$$

This upper bound is always greater than or equal to  $(1 - \frac{1}{m})$ . Thus, it is of limited value in practice, since it says that the mean error rate per labeller is no worse than  $1 - \frac{1}{m}$ , which in turn, for reasonably-sized  $m$ , is quite close to 1 (the trivial upper bound) .

## 4 Application of the Lower Bound

### 4.1 Catalog Generation in Scientific Applications

In a number of observational sciences such as astronomy and planetary geology, a common step in the scientific process is to convert raw data (such as images) into a catalog of objects of interest (Fayyad et al., 1996). Such catalogs form a standard data product which can be used by other scientists as the basis for quantitative scientific studies (such as investigations of the spatial clustering patterns of objects, etc.). Examples include counting stars and galaxies in telescope images to generate a sky catalog, counting impact craters on the surface of the moon, counting and characterizing sunspots in images of the Sun, and counting volcanoes in radar images of Venus. Typically the cataloging is carried out by known experts in the field.

In each of these applications, the quality of the final catalog is inevitably a function of the subjective nature of the cataloging process. In some applications there may be little variation between the labels provided by different experts for the same object: in other applications the variance may be quite high, indicating that the data in the catalog should be treated accordingly. The variation in expert opinion may be due to visual ambiguity introduced by the resolution limits of the data, perhaps the pixel-resolution of an imaging instrument.

### 4.2 Bounding the Mean Accuracy of Volcano Counting

The Magellan spacecraft orbited Venus from 1990 to 1994 and transmitted back to Earth a high resolution synthetic aperture image map of the planet, approximately 30,000 1Mbyte images in total. The study of volcanic features on the surface of Venus is a key issue in planetary geology due to the predominance of volcanism on the planet (Saunders et al., 1992). Generating a comprehensive volcano catalog from the Magellan data is a prerequisite for more advanced studies such as cluster analysis of the volcano locations. Of interest in the context of this paper is the accuracy of the volcano labels provided by planetary geologists.

In previous work a pattern recognition system for automatically counting volcanoes in the Magellan images of Venus has been developed: the pattern recognition system is described in detail elsewhere (Burl et al., 1994a) and is not of direct interest here. As part of the development of the pattern recognition system, several planetary geologists, considered experts in Venus volcanism, provided labels for sets of Venus images as training and test data. Significant variability between the geologist’s labellings was noticed, thus motivating work on the problem of quantifying classification accuracy of both humans and algorithms in the absence of ground truth. The variability in the labelling appears to be primarily due to the relatively low signal-to-noise ratio (relative to small volcano structure) in the SAR images (Fayyad et al., 1996).

Each geologist examined sets of images independently and used mouse-clicks within a graphical user-interface to indicate their estimate of where the volcanoes were located within a given image. The first labelling experiment consisted of 4 images and 4 experts (geologists A, B, C, and D). Between the 4 geologists, 269 estimated volcano locations were found in total in the 4 images. Consider this to be the database of  $N = 269$  items with binary labels: volcano or non-volcano. One can think of each “item” as a local pixel window or region of interest. The lower bound on mean error rate (using Equation (4)) was found to be 19.3%, i.e., the average error rate among geologists A, B, C, and D, labelling volcanoes on these particular 4 Magellan images, is at least 19.3% relative to ground truth .

The second labelling experiment consisted of 2 geologists (A and B from the first experiment) who each individually labelled 38 images (different from the first 4). In this case 512 possible volcano locations were found in total. Again, considering this to be a database of  $N = 512$  items with binary labels results in a lower bound on the mean error rate of A and B of 24.1%. If only the labellings of geologists A and B are considered on the 4 images in the first experiment, they made at least 22.2% errors on average (for these 2 geologists on these 4 images).

Across different subsets of images, with different sets of geologists, the results for the volcano problem have consistently shown a lower-bound on the mean error rate of about 20%. Thus, one can state that at least one of the expert geologists is in error at least 20% of the time in terms of volcano labelling, over a range of different Magellan images, relative to the ground truth. The true mean error rate for the geologists could in fact be much higher than 20%.

### 4.3 Significance of the Results

There are two primary results from applying this approach in general:

1. From a scientific viewpoint, interpretation of subjectively-derived catalogs (such as volcano catalogs derived from the Magellan data set) should take into account the level of disagreement between the labellers. For example, statistics of scientific interest derived from the catalog, such as the mean number of volcanos observed in a given region of the planet, must be interpreted appropriately.
2. From a pattern recognition viewpoint, the bounds may indicate that subjectively-labeled training and test data sets (for any one expert, or consensus of experts) could contain a

significant degree of noise in the labels. Methods for taking this noise into account both in training and evaluation are described in Burl et al. (1994b) and Smyth et al. (1996).

## 5 Comments on Related Work

Previous work on modelling noise in class labels has largely relied on parametric models of the noisy labelling process. For example, *rating models* assume that a set of labellers provide a discrete set of ratings of the likelihood that an item belongs to a class, and from the ratings of multiple labellers an overall combination model and posterior estimates for individual items are found (French, 1985; Agresti, 1992; Uebersax, 1993). A key issue is the nature of the assumptions about the independence of the different labellers.

In another distinct approach, the error patterns in class labelling are assumed to obey a particular model and the implications are analysed (Aitchison and Begg, 1976; Titterton, 1989; Lugosi, 1992). For example, as the noise in the labelling process increases the effect on the estimation performance of certain parametric classification methods has been investigated (Krishnan and Nandy, 1990).

Both of these general approaches bear some relation to the problem discussed in this paper and indeed the ratings approach has been used with success on the volcano data (Smyth et al., 1996). However, the results in this paper are distinct from this prior work in the sense that the bounds derived here make no assumptions whatsoever about the nature of the labelling errors, the independence relationships between the  $K$  labellers, or the underlying distribution of the data.

## 6 Conclusion

Simple bounds were derived for the mean classification error rate of  $K$  labellers in the absence of ground truth. The lower bound was applied to data from a remote-sensing image analysis problem. The results confirmed that the subjective error rate for the problem is quite high. The method has applications to classification problems where data must be labelled in a subjective manner by experts and there is no ground truth available to calibrate their performance.

## Acknowledgements

The author gratefully acknowledges useful discussions on this topic with Michael Burl, Usama Fayyad, Pietro Perona, Jayne Aubele, and Larry Crumpler. The research described in this paper has been carried out by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration.

## References

- Agresti, A. (1992). Modelling patterns of agreement and disagreement. *Statistical Methods in Medical Research*, vol.1, pp.201–218.
- Aitchison, J. and C. B. Begg (1976). Statistical diagnosis when basic cases are not classified with certainty. *Biometrika*, 63 (1), pp.1–12.
- Burl, M. C., U. M. Fayyad, P. Perona, P. Smyth, and M. P. Burl (1994a). Automating the hunt for volcanoes on Venus. In *Proceedings of the 1994 Computer Vision and Pattern Recognition Conference, CVPR-94*, Los Alamitos, CA: IEEE Computer Society Press, pp.302–309.
- Burl, M. C., U. M. Fayyad, P. Perona, and P. Smyth (1994b). Automated analysis of radar imagery of Venus: handling lack of ground truth. In *Proceedings of the IEEE Conference on Image Processing*, Piscataway, NJ: IEEE Press, vol.III, pp.236–240.
- Fayyad, U. M., P. Smyth, M. C. Burl, and P. Perona (1996). Learning to catalog science images. In *Early Visual Learning*, S. Nayar and T. Poggio (eds.), New York, NY: Oxford University Press, pp. 237–268.
- French, S. (1985). Group consensus probability distributions: a critical survey. In *Bayesian Statistics 2*, J. M. Bernardo, M. H. DeGroot, D. V. Lindley, A. F. M. Smith (eds.), Amsterdam, North-Holland: Elsevier Science Publishers, pp.183–202.
- Krishnan, T. and S. C. Nandy (1990). Efficiency of discriminant analysis when initial samples are classified stochastically. *Pattern Recognition*, vol.23, no.5, pp.529–537.
- Lugosi, G (1992). Learning with an unreliable teacher. *Pattern Recognition*, vol. 25, no.1, pp.79–87.
- Saunders, R. S., et al. (1992). Magellan mission summary. *Journal of Geophysical Research*, vol.97, no.E8, pp.13067–13090.
- Smyth, P., M. C. Burl, U. M. Fayyad, and P. Perona (1996). Modeling subjective uncertainty in image annotation. In *Advances in Knowledge Discovery and Data Mining*, U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, R. Uthurasamy (eds.), AAAI/MIT Press, pp.517–539.
- Smyth, P. (1995). Learning with probabilistic supervision. In *Computational Learning Theory and Natural Learning Systems 3*, T. Petsche, S. Hanson, and J. Shavlik, Cambridge, MA: MIT Press, pp.163–182.
- Titterton, D. M. (1989). An alternative stochastic supervisor in discriminant analysis. *Pattern Recognition*, vol.22, no.1, pp.91–95.
- Uebersax, J. S. (1993). Statistical modeling of expert ratings on medical treatment appropriateness. *J. Amer. Statist. Assoc.*, vol.88, no.422, pp.421–427.