

On Methods for Incorporating Evidences into Posterior Scoring of Hypotheses

Dmitry Kovalenko and Sargur N. Srihari
 CEDAR, Department of Computer Science and Engineering
 University at Buffalo, Buffalo NY 14260-2500, USA
 dmitryse@buffalo.edu, srihari@cedar.buffalo.edu

Abstract

In likelihood ratio based decision methods, often a variable number of input evidences is used. A decision based on many such inputs can result in nearly the same likelihood ratio as one based on few inputs. We consider methods for distinguishing between such situations. One of these is to provide confidence intervals together with the decisions and another is to combine the inputs using weights. We suggest a new method that generalizes Bayesian approach and uses an explicitly defined discount function. Empirical evaluation shows greater flexibility of the designed method.

1 Introduction

Consider statistical decision methods based on likelihood ratio. One of the widely used is Bayesian approach to statistical inference of posterior probability of a hypothesis given its prior and a set of evidences (data samples or features). Assume that: h_0 – a hypothesis of interest, h_1 – the alternative hypothesis to h_0 , $F = \{f_i\}_{i=1}^n$ – a set of n observed features. Then:

$$P(h_0|F) + P(h_1|F) = 1 \quad (1)$$

Applying Bayes' rule and assuming that all features are independent we get:

$$P(h_0|F) = \frac{P(h_0) \cdot \prod_i LR(f_i|h_0, h_1)}{P(h_1) + P(h_0) \cdot \prod_i LR(f_i|h_0, h_1)} \quad (2)$$

where: $LR(f_i|h_0, h_1) = \frac{P(f_i|h_0)}{P(f_i|h_1)}$ – is a likelihood ratio of f_i for h_0 to h_1 .

A variable number of features can be used. Nearly the same likelihood ratio value $LR(F|h_0, h_1)$ (and the same overall result) can be based on many features and on few of them. For example, assume the following two cases:

1. A discrete uniform prior ($P(h_0) = P(h_1) = 0.5$) and one feature $F = f_1$ with $LR(f_1|h_0, h_1) = 96$.
2. A discrete uniform prior ($P(h_0) = P(h_1) = 0.5$) and nine features $F = \{f_i\}_{i=1}^9$ with likelihood ratios $\{3, 4, 2, \frac{1}{4}, 2, 2, \frac{1}{3}, 6, 2\}$ accordingly.

For the first case we get: $P_1(h_0|F) = \frac{96}{97}$, which is the same as $P_2(h_0|F) = \frac{96}{97}$ – a posterior probability of h_0 in the second case. Consider an application in Questioned Document Examination (QDE), where: h_0 – two compared documents were written by the same person, h_1 – written by different people. In the first case the result is based on one feature (a discriminating characteristic in QDE) and could be just a chance, whereas in the second case the support of h_0 tends to be more trustworthy since it is based on nine features. We consider methods for distinguishing between such cases. This problem is especially important in case when probabilities of features are estimated approximately.

One of these methods is to use not a single numerical value for posterior probability of h_0 but a credible interval [2]. When there are no parameters of h_0 then we can assume that there is a distribution of likelihood ratios and estimate a credible interval for it. A crucial part of this approach is to define a proper prior probability for likelihood ratios.

Another way of dealing with the problem at hand is to use statistical hypothesis testing [3]: to define likelihood ratio confidence interval with a given level of confidence. If we assume a log-normal distribution with unknown mean and variance as a distribution of likelihood ratios then we can calculate confidence interval for log-likelihood ratio mean and use it in equation 2.

A completely different approach to the problem is to use a score instead of posterior probability of h_0 . This approach allows us not to make any assumptions about likelihood ratio distribution. We have designed such scoring method based on Bayesian approach with

weighted likelihoods and a discount function that diminishes the score when features in a comparison have a small sum of weights. This approach is discussed in this paper in full details.

A set of experiments was conducted on data sets derived from different distributions. Results of these experiments, pros and cons of each described approach are discussed in this paper.

2 Problem Formulation

Assume that there are:

- Two hypotheses: h_0 and its alternative h_1 (e.g. h_0 – two handwritten documents were written by the same author; h_1 – written by different people).
- Prior probability values for hypotheses: $P(h_0)$, $P(h_1)$ (e.g. discrete uniform distribution for two alternative hypotheses: $P(h_0) = P(h_1) = 0.5$).
- Set $F = \{f_i\}_{i=1}^n$ of n extracted features (e.g. F is a set of discriminating characteristics in compared documents).

The problem is to get a score for hypothesis h_0 that reflects how likely h_0 is true in comparison with its alternative and that explicitly includes how trustworthy a hypotheses comparison is.

Here we assume that all features F are independent from each other¹. A common way to solve the problem is to use Bayesian approach (see equation 2). To adapt its ideas in order to explicitly include a hypotheses comparison trustworthiness we discuss several existing methods and suggest a new one.

3 Methods for Relative Scoring of Hypotheses

Assume that likelihood ratios of features are themselves i.i.d. samples of some distribution. For example, we can assume that it is a log-normal distribution. Then it is possible to identify a credible interval [1]: an $\alpha \cdot 100\%$ credible interval is a set C such that $P(C) = \alpha$. There are different ways to evaluate credible intervals approximately or exactly when it is possible [1]. Values of equation 2 determined on credible interval boundaries together with level α define a result for the stated problem.

An analogy in some way of a credible interval in statistical hypothesis testing is a confidence interval [3]. It

¹Adapting discussed in this paper methods to the case when features have certain dependencies is a direction for future research.

is an interval $(L_\alpha(F), R_\alpha(F))$ evaluated from a given set of features F , that frequently includes the parameter of interest, if the experiment is repeated. The frequency here is determined by confidence level α . For example, the parameter of interest could be mean of log-likelihood ratios if they comply with the log-normal distribution. Values of equation 2 determined on confidence interval boundaries along with confidence level α define a solution for the stated problem.

A useful property of using credible or confidence intervals is that in most cases with increase in the number of features the distance between interval boundaries is decreasing. A problem of using credible or confidence intervals is a necessity of combining interval boundaries and level α into one score if only one single value is required by the application (as a score for hypothesis h_0).

To drop any assumptions on likelihood ratio distributions, we have designed a new method to solve the problem. To include credibility of the hypotheses comparison into the score we suggest to use the following expression instead of posterior probabilities ratio:

$$S \left(\frac{h_0|F}{h_1|F} \right) = \frac{P(h_0)}{P(h_1)} \cdot \left(\prod_{i=1}^n \left(\frac{P(f_i|h_0)}{P(f_i|h_1)} \right)^{w_i} \right)^{d(\sum_i w_i)} \quad (3)$$

where: $S \left(\frac{h_0|F}{h_1|F} \right)$ – is the relative score of hypothesis h_0 against h_1 . The following term:

$$\prod_{i=1}^n \left(\frac{P(f_i|h_0)}{P(f_i|h_1)} \right)^{w_i} \quad (4)$$

is a weighted likelihood ratio [5], where weights $w_i \in \mathbb{R}^+$. A trivial set of weights is:

$$w_i = 1, \quad i = \overline{1, n} \quad (5)$$

which defines that all of the used features are equal in their contribution to the overall value. A more discriminating approach is to use higher weights for more accurate and influential features.

Function $d(\sum_i w_i) : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a discount function which reflects a credibility of the comparison. It ought to give a support to a comparison with big sum of features weights and discount the score for a comparison with the small sum. A trivial discount function is:

$$d_0 \left(\sum_i w_i \right) = \frac{n}{\sum_i w_i} \quad (6)$$

which does not make any discount but normalizes weighting scheme. An example of discount function which makes a discount for a small sum of weights is:

$$d_1 \left(\sum_i w_i \right) = \frac{n}{\max(m, \sum_i w_i)} \quad (7)$$

where: m is a constant threshold. When trivial weights are used and a number of features is more than a threshold then score S is equal to the ratio of posterior probabilities defined in equation 2.

Another example of discount function applies different discounts for different values of sum of weights:

$$d_2 \left(\sum_i^n w_i \right) = \beta \cdot \log \left(\sum_i^n w_i \right) \quad (8)$$

where: β – is a magnitude parameter.

Accuracy of using score S instead of posterior probability ratio can be at least the same since the ratio is a special case of S with trivial weights and the trivial discount function. Another advantage of using score S is absence of additional assumptions in comparison with using credible or confidence intervals. The disadvantage however is that in general it results in a score for the hypotheses of interest in comparison with its alternative, but not in the posterior probability.

4 Experimental Results

To reveal pros and cons of the described methods we have conducted a set of experiments. Synthetically created data sets were used since it allows to compare methods on different (initially specified) sample distributions and with a wide range of data sets sizes. Two types of data sets were used:

- I) Data set consists of values derived from normal distribution $\mathcal{N}(50, 25)$ with mean $\mu = 50$ and variance $\sigma^2 = 25$.
- II) Data set consists of values derived from uniform distribution $\mathcal{U}[0, 100]$ defined on interval $[0, 100]$.

Data sets of each type consisted of 5 through 125 samples with step size of 5. Every sample in each data set represents a log-likelihood ratio value for one feature.

Throughout testing we use a discrete uniform prior:

$$P(h_0) = P(h_1) = 0.5$$

Hence, we focus on a score for log-likelihood ratio (LLR) rather than on a posterior score for the hypotheses h_0 . The following methods were used in the comparison:

- I) Original Bayesian approach where score is defined as posterior probabilities ratio of the hypotheses.
- II) An approach based on estimation of confidence interval for LLR distribution mean. We assume that likelihood ratios have a normal distribution $\mathcal{N}(\mu, \sigma^2)$ with unknown mean and variance.

- III) Discounted weighted likelihood ratio with trivial weights (see equation 5) and discount function d_1 (see equation 7) and threshold $m = 100$.

We did not include an approach based on credible intervals since if we assume that log-likelihood has a normal distribution with unknown mean and precision and use a normal-gamma distribution as a conjugate prior then (with certain parameters) credible and confidence intervals coincide[4].

Results on data sets generated from normal distribution $\mathcal{N}(50, 25)$ are shown in figures 1(a) and 1(b). Results on data sets generated from uniform distribution $\mathcal{U}[0, 100]$ are shown in figures 1(c) and 1(d). Blue circles show log-likelihood ratio (method I); green solid lines and blue circles on figures 1(b) and 1(d) show 95% confidence intervals and log-likelihood ratio average accordingly (method II); red dots show discounted weighted log-likelihood ratio score (method III).

In both sets of tests method III converges to log-likelihood ratio for large numbers of samples and gives a discounted value for small number of them. In particular, the score average is close to zero for less than 10 samples while LLR average is the same as for large number of samples (see figures 1(b) and 1(d)).

Let us refer to the example from the introduction: comparison of two cases (in the first case one feature is given and in the second - nine features). Applying Bayesian approach we get:

$$P_1(h_0|F) = P_2(h_0|F) = 0.99$$

where: $P_1(h_0|F)$ – posterior probability of h_0 in the first case, which is the same as posterior probability $P_2(h_0|F)$ of h_0 in the second case. Using discounted weighted likelihood ratio with same parameters as earlier (trivial weights and discount function d_1 but with threshold $m = 20$) we get the following posterior probability scores for the first and the second cases:

$$P_1^s(h_0|F) = 0.89, P_2^s(h_0|F) = 0.56$$

The use of the proposed score enables distinguishing between such cases.

Experimental results show that there are the following disadvantages of using confidence intervals: i) a necessity to make an assumption about log-likelihood ratios distribution; ii) confidence interval boundaries are wide (even for 95% confidence level) and especially in case when real data distribution is different from the assumed one (see figure 1(c-d)); iii) there is no common way to combine confidence interval boundaries and confidence level into one score (see figures 1(a),1(c): if we use confidence interval boundaries times the number of samples then the result scores become even more distant).

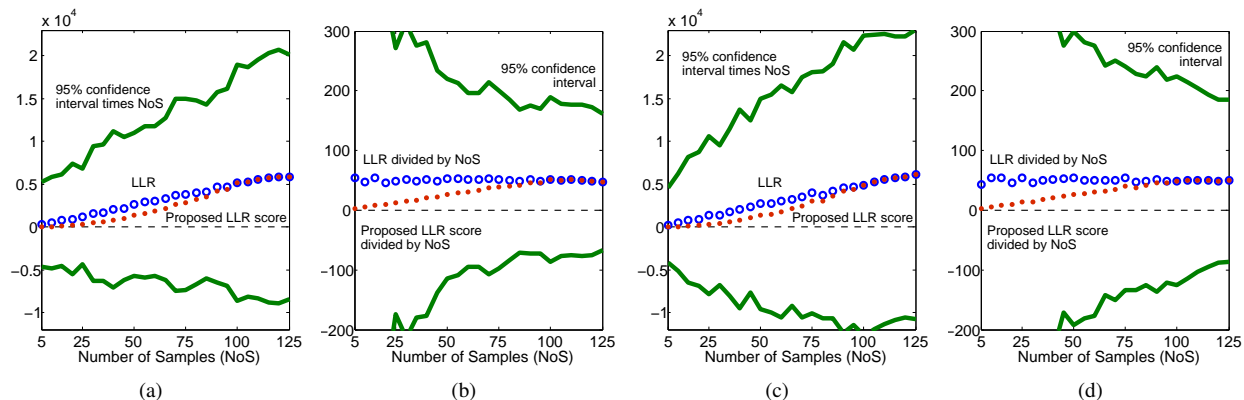


Figure 1. Scores on normal distribution $\mathcal{N}(50, 25)$ (a-b) and uniform distribution $\mathcal{U}[0, 100]$ (c-d): i) Green solid lines – 95% confidence intervals for the log-likelihood ratio mean; ii) Blue circles – log-likelihood ratio; iii) Red dots – discounted weighted log-likelihood ratio score with trivial weights (eq. 5) and discount function d_1 (eq. 7) and threshold $m = 100$.

The advantages of using a discounted weighted likelihood are: i) it generalizes Bayesian approach and allows to explicitly set a discount for hypotheses comparison based on a few features; ii) no need to make additional assumptions on likelihood ratios distribution; iii) flexibility of choosing weights and discount function allows tuning the approach for every particular application. However, there are disadvantages also: i) the likelihood score results not in a posterior probability but in a relative score; ii) the flexibility of choosing weights and a discount function makes it non-trivial to find the best settings for a particular application.

5 Summary and Conclusion

The problem of estimating a posterior score of a hypothesis given prior beliefs and a set of evidences (observed features) has been considered. A common way of using posterior probability as the score has a drawback. Nearly the same result can be based on many features and on few of them, which can be important to distinguish. We consider methods for distinguishing between such situations. One of these is to provide confidence intervals for the obtained result and another is to combine features using weights and a discount function.

The usage of intervals has the following problem: there is no common way to combine interval boundaries and the corresponding level into one score.

The designed method has the following advantages: i) flexibility (it can be tuned for every particular application); ii) accuracy of the designed approach at least the same as of Bayesian approach; iii) there are no assumptions on likelihood distribution (in comparison to

approaches based on credible or confidence intervals). Although the scoring method does not in general provide probabilities as output (which is a disadvantage), it distinguishes the relative strengths of having different numbers of features.

Automated methods for selecting a weighting scheme and choosing a discount function are topics for future research.

Acknowledgement

The work was supported in part by Award No. 2010-DN-BX-K037 of the NIJ, Office of Justice Programs, U.S. DoJ. The opinions expressed are those of the author(s) and do not necessarily reflect those of the DoJ.

References

- [1] G. Antelman, A. Madansky, and R. McCulloch. *Elementary Bayesian Statistics*. Cheltenham, UK, 1997.
- [2] P. Congdon. *Bayesian Statistical Modelling*. Wiley, USA, 2007.
- [3] D. Cox and D. Hinkley. *Theoretical Statistics*. Chapman and Hall/CRC, USA, 1979.
- [4] K. Murphy. *Conjugate Bayesian analysis of the Gaussian distribution*. UBC, 2007.
- [5] X. Wang. Asymptotic properties of maximum weighted likelihood estimators. *Journal of Statistical Planning and Inference*, 119(1):37–54, 2004.