

One-Class Multiple-Look Fusion: a theoretical comparison of different approaches with examples from infrared video

Mark W. Koch

Sandia National Laboratories[†]
Albuquerque, NM 87185-1163
mwkoch@sandia.gov

Abstract

Multiple-look fusion is quickly becoming more important in statistical pattern recognition. With increased computing power and memory one can make many measurements on an object of interest using, for example, video imagery or radar. By obtaining more views of an object, a system can make decisions with lower missed detection and false alarm errors. There are many approaches for combining information from multiple looks and we mathematically compare and contrast the sequential probability ratio test, Bayesian fusion, and Dempster-Shafer theory of evidence. Using a consistent probabilistic framework we demonstrate the differences and similarities between the approaches and show results for an application in infrared video classification.

1. Introduction

There have long been multiple competing approaches for accomplishing multiple look sensor fusion. By multiple look fusion we assume we can make multiple measurements on an object, as it passes through the field of view of the sensor. For example, multiple frames in an infrared video or the extraction of multiple high-resolution-range profiles from ground-moving-target-indicator radar can produce multiple sensor measurements on an object.

Here, we will compare and contrast three data fusion approaches: sequential probability ratio test (SPRT), Bayesian, and Dempster Shafer (DS). To accomplish this we will use a common probabilistic framework that is useful in real-world pattern recognition problems in unconstrained environments. Our goal is not to say one is better than the other, but to find commonalities between the various approaches and use one approach to find

insights into others. This objective leads to a more unified approach in sensor fusion that can take advantage of the best features for all the approaches.

2. One Class Classifiers

A popular assumption in most pattern recognition and multilook fusion approaches is the “closed-world” assumption [3]. Here, all the objects or actions sequences are assumed to be known. This allows the use of “key features” to distinguish objects/actions from each other and early termination of a sequential test when one of these key features is detected [3]. Unfortunately, this approach does not work in unconstrained environments where potentially any moving object can appear in a scene.

For unconstrained environments, a structure one might choose for a classification problem is based on a *one-class* classifier [1]. For one class classifiers, we are interested in one specific target θ_1 represented by the alternative hypothesis H_1 , and the null hypothesis H_0 represents the nontarget $\bar{\theta}_1$ class. While this might seem like an oversimplification, one could argue that for a multi-class problem with other targets of interest $\theta_2, \dots, \theta_m$ one would design a one-class classifier for each of them. For the θ_1 one-class classifier we can further divide the nontargets into two groups: the other targets of interest $\theta_2, \dots, \theta_m$ and the unknown class θ_0 . This allows us to further distinguish between two types of false alarm errors: *between-class* and *out-of-class*. Between-class errors occur when alarming on another target $\theta_2, \dots, \theta_m$ by calling it the target θ_1 . Out-of-class errors occur when alarming on an unknown signature θ_0 by calling it the target θ_1 . In making any decision, we want to control two types of errors: *missed detection errors* and *false alarm errors*. Missed detection errors result from missing a target signature by calling it a nontarget, and false alarm errors result from alarming on a nontarget signature by calling it a target.

For example, suppose we are interested in classifying moving objects in infrared video as humans, vehicles, or unknown. Here, θ_1 would represent the human class and

[†] Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy’s National Nuclear Security Administration under Contract DE-AC04-94AL85000.

θ_2 the vehicle class. The unknown class θ_0 would represent all moving objects not in θ_1 or θ_2 . This could be wind-blown objects (tumbleweed, boxes, trash cans, etc.) or animals. These unknown moving objects, a possible source of the out-of-class errors, are a significant problem in real-world pattern recognition problems in unconstrained environments. A Bayesian classifier approach designed for the human vs. vehicle problem, while minimizing the between-class errors, would require models of all the possible objects that could be imaged by the sensor to control the out-of-class errors. Otherwise it would classify an animal as human or vehicle. Modeling “the whole world” of possible objects is untenable for most realistic systems deployed in unconstrained environments. Instead, we use a goodness of fit (GOF) classifier to control the out-of-class errors, and power analysis [9] to model the unknown class.

Whereas Bayesian classifiers minimize the between-class error, they do nothing to control the out-of-class errors. Figure 1 illustrates this potential problem. The figure shows a two-dimensional feature space, with samples from two targets: Target A represented by stars and Target B represented by filled circles. Assuming normal distributions and equal covariance matrices for the targets, the Bayes decision boundary has a linear form (Figure 1a). Whereas the Bayes classifier minimizes the between-class errors of the A and the B targets, it does not control the out-of-class errors caused by unknown objects represented by “x” symbols. Depending on which side of the boundary the nontarget falls, the classifier will assign the unknown to one of the known classes and make an out-of-class error. Figure 1b shows a GOF classifier that tries to surround the target class. Here, the unknown objects, that have widely differing features from the target class (“x” symbols), will be classified correctly. In general, the GOF classifier have improved out-of-class errors, but the between-class errors will increase, since it is not necessarily an optimal Bayes classifier.

3. Probabilistic Framework

We will use the same probabilistic framework for comparing and contrasting the fusion algorithms. We start by assuming a stream of observations represented as random variables x_1, x_2, \dots . These samples result from the best match scores in the GOF metric. We also assume knowledge of the functions $p(x_i | T)$ and $p(x_i | \bar{T})$ which represent the probability density function (PDF) of an observation x_i , given the target T and the nontarget \bar{T} , respectively. The PDF’s can be discrete or continuous and can be determined through theoretical means or empirical modeling.

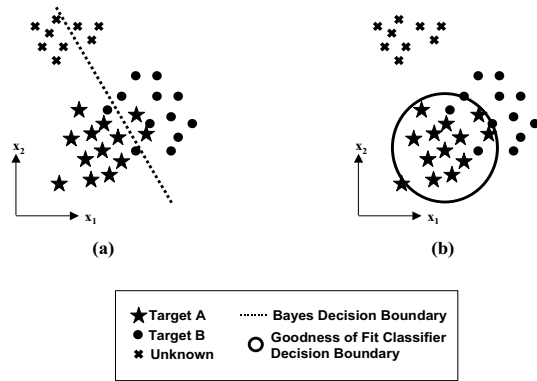


Figure 1: Comparison of Bayes and goodness of fit (GOF) classifiers. (a) Bayes classifier. (b) GOF classifier.

The PDF $p(x_i | T)$ is usually straightforward to determine, since one has knowledge of the target of interest. This information can come from data collections or modeling and simulation using CAD descriptions and a physics-based sensor signature prediction software such as Irma [10]. The PDF $p(x_i | \bar{T})$ is usually more problematic. One approach for modeling the nontarget class uses *statistical power analysis* [9] to model the *worst case nontarget*. The approach has some similarities to that taken by [1] for modeling composite hypotheses by determining the *least favorable choice*. Power analysis assumes the tested effect is linear and the measured effect size (small, medium or large) is known. Typically, power analysis allows the statistician to determine if enough samples were collected to give the test a high power. While the exact form depends on $p(x_i | T)$, we show an example from [6] where $p(x_i | T)$ is $N(0,1)$. This is convenient in problems where the central limit theorem can be applied. In [6], it was shown that the worst case nontarget PDF is $N(\mu_N, 1)$. Here the location parameter μ_N represents the smallest acceptable effective difference between the target and nontarget.

4. Sequential Probability Ratio Test

One approach for combining match scores as they become available is to use the Wald sequential hypothesis test or SPRT [12]. After n observations the likelihood ratio is:

$$\Lambda(n) = \prod_{i=1}^n \lambda_i, \text{ where } \lambda_i = \frac{f(x_i | T)}{f(x_i | \bar{T})} \quad (1)$$

Often it is more numerically convenient to work with the log-likelihood:

$$Z(n) = \log(\Lambda(n)) = \sum_{i=1}^n z_i, \text{ where } z_i = \log\left(\frac{p(x_i | T)}{p(x_i | \bar{T})}\right) \quad (2)$$

The SPRT uses two decision boundaries (a, b) to make a decision:

$$\begin{aligned} \text{Reject } H_0 & \quad \text{If } Z(n) \geq a \\ \text{Accept } H_0 & \quad \text{If } Z(n) \leq b \\ \text{Get more data} & \quad \text{If } b < Z(n) < a \end{aligned} \quad (3)$$

One desirable property of the SPRT is that the decision boundaries can be determined from the desired error rates. Thus, these decision boundaries can be obtained using the desired false alarm rate, α , and the desired missed detection rate, β :

$$a = \log\frac{1-\beta}{\alpha} \text{ and } b = \log\frac{\beta}{1-\alpha}. \quad (4)$$

It has been shown that the SPRT, on average, uses the smallest number of observations to make a decision [12].

5. Bayesian Fusion

A sequential update formula can be derived from Bayes formula:

$$p(T | x_1) = \frac{p(T)p(x_1 | T)}{p(T)p(x_1 | T) + (1-p(T))p(x_1 | \bar{T})} \quad (5)$$

The quantity $p(T | x_1)$ is the posterior probability of a target after one observation x_1 . By substituting $p(T | x_1)$ for the prior $p(T)$ and a new observation x_2 for x_1 in (5) we get the posterior probability of a target after two observations as:

$$\frac{p(T)p(x_1 | T)p(x_2 | T)}{p(T)p(x_1 | T)p(x_2 | T) + (1-p(T))p(x_1 | \bar{T})p(x_2 | \bar{T})} \quad (6)$$

Using $Y(n)$ to represent the Bayesian posterior probability after n observations or $p(T | x_1 \dots x_n)$ we get

$$Y(n) = \frac{\lambda_0 \Lambda(n)}{\lambda_0 \Lambda(n) + 1} \quad (7)$$

where $\Lambda(n)$ is the Wald likelihood ratio in (1) and λ_0 represents the a-priori likelihood ratio:

$$\lambda_0 = \frac{p(T)}{(1-p(T))}. \quad (8)$$

The Bayesian posterior probability is always between 0 and 1. By (1) and the definition of a PDF $0 \leq \Lambda(n) < \infty$, so as $\Lambda(n) \rightarrow \infty$ then $Y(n) \rightarrow 1$ and if $\Lambda(n) = 0$ then $Y(n) = 0$. We can use the SPRT stopping conditions to determine thresholds on the Bayesian posterior probability. If

$$a = \log(A) \text{ and } b = \log(B) \quad (9)$$

then A and B are the upper and lower stopping conditions for $\Lambda(n)$, respectively. Thus the Bayesian stopping rule becomes:

$$\begin{aligned} \text{Reject } H_0 & \quad \text{If } Y(n) \geq C \\ \text{Accept } H_0 & \quad \text{If } Y(n) \leq D \\ \text{Get more data} & \quad \text{If } C < Y(n) < D \end{aligned} \quad (10)$$

where

$$C = \frac{\lambda_0}{1+\lambda_0}, A = \frac{\lambda_0}{1+\lambda_0} \frac{(1-\beta)}{\alpha} \quad (11)$$

and

$$D = \frac{\lambda_0}{1+\lambda_0}, B = \frac{\lambda_0}{1+\lambda_0} \frac{\beta}{1-\alpha}. \quad (12)$$

Note the $\lambda_0/(1+\lambda_0)$ basically tweaks the threshold according to the a-priori information. For $\lambda_0 > 1$ the thresholds will go lower to make target calls more probable and if $\lambda_0 < 1$ the thresholds will go higher to make nontarget calls more probable.

6. Dempster-Shafer Theory of Evidence

The Dempster Shafer (DS) theory [11] is a mathematical theory of evidence that allows combining evidence from different sources to arrive at a degree of belief. It models uncertainty by not requiring one to assign all of one's belief to a single proposition.

The main assumption we make is that evidence is *consonant*. This allows us to use the probabilistic framework that we established in section 3. Shafer defines consonant evidence as evidence that points in a single direction and only varies in the precision of focus [11]. This fits well with the GOF metric. The GOF describes the difference between stored knowledge, for example a template of the target, and the measured data. Thus it points only in the direction and focus of the hypothesis represented by the stored knowledge.

For the one-class problem the frame of discernment or set of possible outcomes is $\{\phi, T, \bar{T}, \Theta\}$ where $\Theta = \{T, \bar{T}\}$ and ϕ represents the empty set. From [11] the support function for the target T is

$$\left. \begin{aligned} m_x(\bar{T}) &= 0 \\ m_x(T) &= 1 - \frac{p(x|\bar{T})}{p(x|T)} \\ m_x(\Theta) &= \frac{p(x|\bar{T})}{p(x|T)} \end{aligned} \right\} \text{if } \frac{p(x|T)}{p(x|\bar{T})} \geq 1 \quad (13)$$

and the support for the nontarget \bar{T} is

$$\left. \begin{aligned} m_x(\bar{T}) &= 1 - \frac{p(x|T)}{p(x|\bar{T})} \\ m_x(T) &= 0 \\ m_x(\Theta) &= \frac{p(x|T)}{p(x|\bar{T})} \end{aligned} \right\} \text{if } \frac{p(x|\bar{T})}{p(x|T)} > 1. \quad (14)$$

where $m(\cdot)$ represents the DS basic probability assignment (BPA) function and $m_x(\Theta)$ represents the amount of uncertainty in the observation x .

For what we want to show it is simpler to work with DS's weight of evidence. If $m(A)$ represents the BPA for the proposition $A \subseteq \Theta$ then weight of evidence $w(A)$ is:

$$w(A) = -\log(1 - m(A)). \quad (15)$$

In terms of weight, of evidence equation (13) becomes:

$$\left. \begin{aligned} w_x(\bar{T}) &= 0 \\ w_x(T) &= \log\left(\frac{f(x|T)}{f(x|\bar{T})}\right) \end{aligned} \right\} \text{if } \frac{f(x|T)}{f(x|\bar{T})} \geq 1 \quad (16)$$

and equation (14) becomes

$$\left. \begin{aligned} w_x(\bar{T}) &= \log\left(\frac{f(x|\bar{T})}{f(x|T)}\right) \\ w_x(T) &= 0 \end{aligned} \right\} \text{if } \frac{f(x|\bar{T})}{f(x|T)} > 1 \quad (17)$$

Using equation (2)

$$\left. \begin{aligned} w_x(\bar{T}) &= 0 \\ w_x(T) &= z \end{aligned} \right\} \text{if } \frac{f(x|T)}{f(x|\bar{T})} \geq 1 \quad (18)$$

and

$$\left. \begin{aligned} w_x(\bar{T}) &= -z \\ w_x(T) &= 0 \end{aligned} \right\} \text{if } \frac{f(x|\bar{T})}{f(x|T)} > 1 \quad (19)$$

where z represents the log-likelihood ratio in the SPRT. As long as we separate the evidence and combine only evidence supporting the same proposition, then we have the homogenous weight of evidence combination rule where the weights of evidence combine additively. For the total support of the target class T , let w^+ represent the total amount of positive weight of evidence. Similarly define w^- for the total support for the nontarget class \bar{T} . Here

$$w^+ = \sum_{i=1}^n w_{x_i}(T) = \sum_{i=1}^n z_i, \text{ for } \frac{f(x|T)}{f(x|\bar{T})} > 1 \text{ or } z_i > 0 \quad (20)$$

and

$$w^- = \sum_{i=1}^n w_{x_i}(\bar{T}) = \sum_{i=1}^n z_i, \text{ for } \frac{f(x|\bar{T})}{f(x|T)} > 1 \text{ or } z_i < 0 \quad (21)$$

Combining conflicting weights (w^+ and w^-) of evidence cannot be done by simple addition. From [11] the weight of evidence for the contradictory propositions T and \bar{T} becomes:

$$w(T) = \log\left(\frac{e^{w^+} + e^{w^-} - 1}{e^{w^-}}\right) \quad (22)$$

and

$$w(\bar{T}) = \log\left(\frac{e^{w^+} + e^{w^-} - 1}{e^{w^+}}\right). \quad (23)$$

Using (15) in reverse we can get the corresponding BPA for propositions T and \bar{T}

$$m(T) = \frac{e^{w^+} - 1}{e^{w^+} + e^{w^-} - 1} \quad (24)$$

and

$$m(\bar{T}) = \frac{e^{w^-} - 1}{e^{w^+} + e^{w^-} - 1} \quad (25)$$

The DS uncertainty is the BPA assigned to Θ or

$$m(\Theta) = \frac{1}{e^{w^+} + e^{w^-} - 1} \quad (26)$$

Since $m(T) + m(\bar{T}) + m(\Theta) = 1$ for a one-class classifier. Note the uncertainty of the system's belief is driven down to 0 as evidence is collected or as w^+ and/or w^- increase. Thus high but equal w^+ and w^- would give low uncertainty, but an uninformed decision.

The most obvious decision rule is

$$\begin{aligned} &\text{Decide } T \quad \text{If } w(T) - w(\bar{T}) \geq 0 \\ &\text{Decide } \bar{T} \quad \text{otherwise} \end{aligned} \quad (27)$$

For the one-class problem, this turns out to be equivalent to the three unambiguous decision rules proposed by Kim in [5]. After some algebraic manipulation one can show that (27) is equivalent to

$$\begin{aligned} &\text{Decide } T \quad \text{If } w^+ - w^- \geq 0 \\ &\text{Decide } \bar{T} \quad \text{otherwise} \end{aligned} \quad (28)$$

This is equivalent to a forced SPRT decision if one is unwilling or unable to wait for any more observations.

7. Results

We show results on a one-class problem for video motion classification (VMC). The target class T represents upright-human dismounts and the nontarget class \bar{T} is any other mover detected by the video motion detection (VMD) algorithm. The imager is an uncooled DRS E3500 infrared camera with 320×240 resolution and 8-bit precision. We have collected over 150 video clips of humans walking, running, and crawling. VMD was accomplished with background subtraction [1] and tracking with an alpha-beta tracker. The features for VMC were histograms of oriented gradients (HOG) [2] and the GOF metric was based on multinomial pattern matching (MPM) [7].

Figure 2 shows a mosaic of chips (subimages containing the detection) collected of a runner at a range of about 135 meters. The chips are of size 29×24 , which is very small, giving on the order of 100 pixels on target. Performance on only one observation per frame is mediocre, especially when the runner is obscured by a brighter object as is the case in the first row of chips in the mosaic.

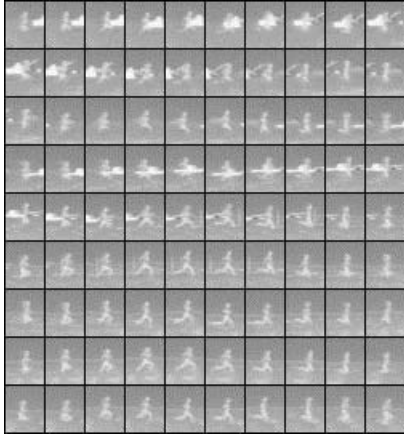


Figure 2: Infrared detections of a runner tracked in an uncooled infrared video imager.

Figure 3 shows the result of multilook fusion using the SPRT, Bayesian Fusion, and DS on the outputs of the MPM GOF classifier. Since MPM is designed to produce $N(0,1)$ scores for HOG features from a target, $p(x_i | T)$ is set to $N(0,1)$. As discussed in Section 3, we use $N(\mu_N, 1)$ for $p(x_i | \bar{T})$ where μ_N is empirically set to 5. The desired error rates α and β are set to 1×10^{-3} .

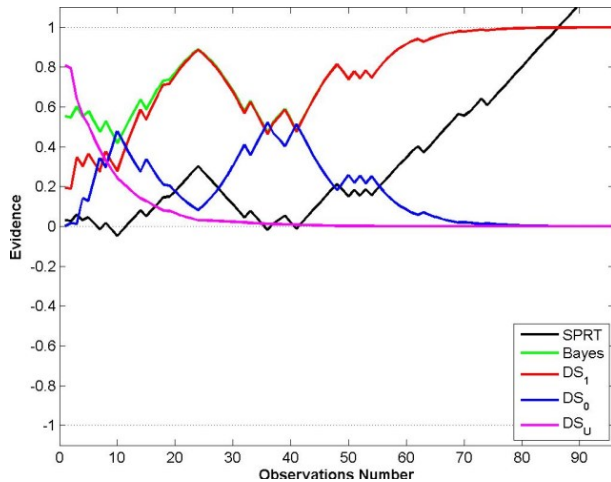


Figure 3: Results of different multilook fusion algorithms of the runner GOF scores from runner chip sequence.

The black line in Figure 3 shows the SPRT cumulative log likelihood $Z(n)$ (2) normalized by the upper SPRT threshold a (4). This puts the SPRT on the same scale as the Bayesian posterior probability and DS BPA and gives a target or nontarget declaration when $Z(n)/a$ passes 1 or -1, respectively (since $a = |b|$). The green line shows the Bayesian posterior probability $Y(n)$ (7) with λ_0 set to 1 (equal priors). The DS BPA is shown by three curves. The red curve is the BPA for the T class $m(T)$, the blue curve is $m(\bar{T})$, and the magenta curve represents the uncertainty of the belief $m(\Theta)$.

Note that as the uncertainty $m(\Theta)$ goes to zero the Bayesian posterior probability $Y(n)$ approaches $m(T)$. Also when the current SPRT Bayesian posterior probability points to a nontarget or $Z(n) < 0$ then $m(T) < m(\bar{T})$ and then for $Z(n) > 0$ $m(T) > m(\bar{T})$. This supports the result in (28) that SPRT and DS make the same forced decision. There is also a similar relation between the Bayes posterior probability $Y(n)$ and $Z(n)$. When $Z(n) < 0$, $Y(n) < 0.5$, and $Z(n) > 0$, $Y(n) > 0.5$. This becomes evident from (7) when $\lambda_0 = 1$ and the knowledge that $Z(n) = 0$ corresponds to $\Lambda(n) = 1$.

It is interesting to note that between observation 30 and 40 the DS uncertainty $m(\Theta)$ is close to 0, but so is the SPRT log-likelihood. Any decision at this point would be an uninformed decision at low uncertainty. When $Z(n)/a$ goes above the threshold of 1, then we can make a high confidence decision that corresponds to low error rates of α and β set to 1×10^{-3} . Also note the Bayesian posterior probability of $Y(n)$ and D.S. BPA $m(T)$ of a target are very close to 1.

Figure 4 shows a mosaic of a chicken tracked by the system. The chip size is 51×34 . Even though the chip size is different than that of the runner the use of the HOG features with the same number of horizontal and vertical blocks gives a system that is scale invariant.

Figure 5 shows the multilook fusion results for the tracked chicken against the upright-human dismounts classification system. When the SPRT log likelihood $Z(n)/a$ falls below the threshold of -1, we can make a nontarget decision with low error rates. Also $Y(n)$ and $m(\bar{T})$ go to zero indicating a nontarget.

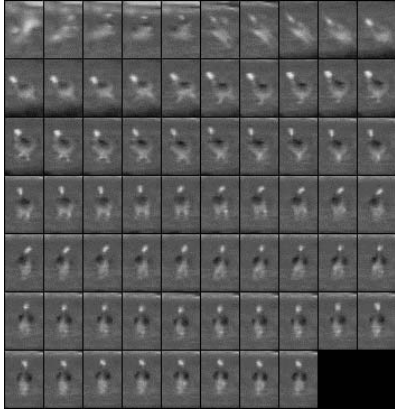


Figure 4: Infrared detections of a chicken tracked in an uncooled infrared video imager.

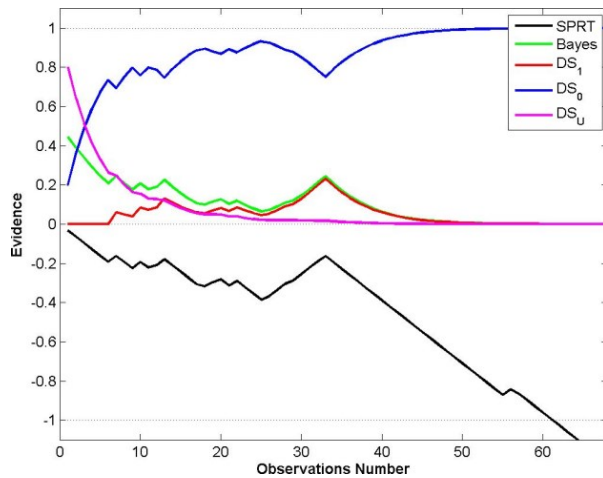


Figure 5: Results of different multilook fusion algorithms of the chicken GOF scores.

8. Conclusion

Here we compared and contrasted three approaches for multilook sensor fusion: the sequential probability ratio test (SPRT), Bayesian fusion, and Dempster Shafer (DS) theory of evidence. The comparison was done in the framework of a one-class classifier using a goodness of fit metric. While a one-class classifier approach may seem severely limiting one could solve a multi-class problem by designing a one-class classifier for each class.

The SPRT is attractive because one can specify the desired missed detection and false alarm errors of a decision. Here, the SPRT remains undecided until enough observations are gathered and the log-likelihood ratio surpasses the a or b thresholds. On the other hand computing a Bayesian posterior probability is attractive

because that is a concept most people can relate to since probabilities are commonly used in weather forecasting. In the paper we showed how to compute SPRT type thresholds, but based on the Bayesian posterior probability. The DS approach is unique in that it incorporates the uncertainty of a belief in its belief combination rule. We showed that for consonant evidence and using the one-class classifier framework this uncertainty decreases as the weight of evidence for either proposition increases or as the number of observations increases. Once this uncertainty approaches zero the DS BPA for a target approaches the Bayesian posterior probability. We also showed that to make a decision with the DS is equivalent to a forced SPRT decision.

References

- [1] E. Boulton, X. Gao, R. Micheals, and M. Eckmann, "Omni-directional visual surveillance," *Image and Vision Computing*, vol. 22, pp. 515, 2004.
- [2] N. Dalal, and B. Triggs, "Histogram of oriented gradients for human detection," *Computer Vision and Pattern Recognition Conference*, June 2005.
- [3] J. W. Davis and A. Tyagi, "Minimal-latency human action recognition using reliable-inference," *Image and Vision Computing*, vol. 24, pp. 455-472, 2006.
- [4] J. D. Gibson and J. L. Melsa, *Introduction to Nonparametric Detection with Applications*. New York: IEEE Press, pp. 25, 1996.
- [5] H. Kim, and P. H. Swain, "Evidential reasoning approach to multisource-data classification in remote sensing," *IEEE Transactions on Systems, Man, and Cybernetics*, **24**, No. 8, 1257-1265, 1995.
- [6] M. W. Koch, G. B. Haschke, and K. T. Malone, "Classifying acoustic signatures using the sequential probability ratio test," *Sequential Analysis Journal*, vol. 23, 4, pp. 557-583, 2004.
- [7] M. L. Koudelka, J. A. Richards, and M. W. Koch, "Multinomial pattern matching for high range resolution radar profiles," *Algorithms for Synthetic Aperture Radar Imagery XIV*. Edited by Zelnio, Edmund G.; Garber, Frederick D.. Proceedings of the SPIE, Volume 6568, 65680V (2007).
- [8] M. Moya and D. Hush, "Network constraints and multi-objective optimization for one-class classification," *Neural Networks*, vol. 9, 3, pp. 463-474, 1996.
- [9] K. Murphy, and B. Myors, *Statistical Power Analysis: A Simple and General Model for Traditional and Modern Hypothesis Tests*. New Jersey: Lawrence Erlbaum Associates, 1998.
- [10] J. Savage, et. al., "Irma 5.2 multi-sensor signature prediction model," *Proc. SPIE 6564, Modeling and Simulation for Military Operations II*, April 2007.
- [11] Shafer, G., *A Mathematical Theory of Evidence*, Princeton University Press, 1976.
- [12] A. Wald, *Sequential Analysis*, John Wiley & Sons Inc, New York, 1947.