

Collaborative Signal Processing for Distributed Classification in Sensor Networks

Ashwin D'Costa and Akbar M. Sayeed*

Electrical and Computer Engineering
University of Wisconsin-Madison
dcosta@cae.wisc.edu, akbar@engr.wisc.edu

Abstract. Sensor networks provide virtual snapshots of the physical world via distributed wireless nodes that can sense in different modalities, such as acoustic and seismic. Classification of objects moving through the sensor field is an important application that requires collaborative signal processing (CSP) between nodes. Given the limited resources of nodes, a key constraint is to exchange the least amount of information between them to achieve desired performance. Two main forms of CSP are possible. Data fusion – exchange of low dimensional feature vectors – is needed between correlated nodes, in general, for optimal performance. Decision fusion – exchange of likelihood values – is sufficient between independent nodes. Decision fusion is generally preferable due to its lower communication and computational burden. We study CSP of multiple node measurements for classification, each measurement modeled as a Gaussian (target) signal vector corrupted by additive white Gaussian noise. The measurements are partitioned into groups. The signal components within each group are perfectly correlated whereas they vary independently between groups. Three classifiers are compared: the optimal maximum-likelihood classifier, a data-averaging classifier that treats all measurements as correlated, and a decision-fusion classifier that treats them all as independent. Analytical and numerical results based on real data are provided to compare the performance of the three CSP classifiers. Entropy comparison between data- and decision-fusion is also provided to quantify the lower communication burden in decision fusion. Our results indicate that the sub-optimal decision fusion classifier, that is most attractive in the context of sensor networks, is also a robust choice from a decision-theoretic viewpoint.

1 Introduction

Wireless sensor networks are an emerging technology for monitoring the physical world with a densely distributed network of wireless nodes [1]. Each node has limited communication and computation ability and can sense the environment in a variety of modalities, such as acoustic, seismic, and infra red [1, 2, 3]. A wide variety of applications are being envisioned for sensor networks, including disaster

* This work was supported by DARPA SensIT program under Grant F30602-00-2-0555.

relief, border monitoring, condition-based machine monitoring, and surveillance in battlefield scenarios. Detection and classification of objects moving through the sensor field is an important task in many envisioned applications. Exchange of sensor information between different nodes in the vicinity of the object is necessary for reliable execution of such tasks due to a variety of reasons, including limited (local) information gathered by each node, variability in operating conditions, and node failure. Consequently, development of theory and methods for collaborative signal processing (CSP) of the data collected by different nodes is a key research area for realizing the vision of sensor networks.

The CSP algorithms have to be developed under the constraints imposed by the limited communication and computational abilities of the nodes as well as their finite battery life. A key goal of CSP algorithms in sensor networks is to exchange the least amount of data between nodes to attain a desired level of performance. In this paper, with the above goal in mind, we investigate CSP algorithms for single-target classification based on multiple acoustic measurements at different nodes. The numerical results presented here are based on real data collected in the DARPA SensIT program.

Some form of region-based processing is attractive in sensor networks in order to facilitate CSP between nodes and also for efficient routing of information in applications involving tracking of moving targets [3]. Typically, the nodes in the network are partitioned into a number of regions and a manager node is designated within each region to facilitate CSP between the nodes in the region and for communication of information from one region to another. Single target classification and tracking generally involves the following steps [3]:

1. **Target detection and data collection.** A target is detected in a particular region which becomes the active region. The detection of a target itself may involve CSP between nodes. For example, outputs of energy detectors may be communicated to the manager node to make the final decision. The nodes within the region that detect the target also collect time series data in different modalities that is communicated to the manager node for classification purposes.
2. **Target localization.** Target detection information (for example, the time of closest point of approach and energy detector outputs) from different nodes is used by the manager node to estimate the location of the target.
3. **Target location prediction.** Location estimates over a period of time are used by the manager node to predict target location at future time instants.
4. **Creation of new potential active regions.** When the target gets close to exiting the current region, the estimates of predicted target location are used to put new regions on alert for target detection.
5. **Determination of new active region.** Once the target is detected in a new region it becomes the new active region. The above four steps are repeated for target tracking through the sensor field.

In this paper, we are primarily concerned with CSP techniques for combining the data collected by different nodes for single-target classification within a par-

ticular active region. However, the basic principles apply to distributed decision making in sensor networks in general.

There are two main forms of information exchange between nodes dictated by the statistics of measured signals. If two nodes yield correlated measurements, *data fusion* is needed, in general, for optimal performance – exchange of (low-dimensional) feature vectors that yield sufficient information for desired classification performance. On the other hand, if two nodes yield independent measurements, *decision fusion* is sufficient – exchange of likelihood values (scalars) computed from individual measurements. In general, the measurements would exhibit a mixture of correlated and independent components and would require a combination of data and decision fusion between nodes. In the context of sensor networks, decision fusion is clearly the more attractive choice. First, it imposes a significantly lower communication burden on the network, compared to data fusion, since only scalars are transmitted to the manager node [3]. Second, it also imposes a lower computational burden compared to data fusion since lower dimensional data has to be jointly processed at the manager node.

In this paper, we investigate the design of CSP classifiers and assess their performance in an idealized abstraction of measurements from multiple nodes. We consider $K = Gn_G$ measurements corresponding to a particular event. The K measurements are split into G groups with n_G measurements in each group. The signal component in the n_G measurements in a particular group is identical (perfectly correlated), but it varies independently from group to group. We compare the performance of three classifiers: 1) the optimal maximum-likelihood (ML) classifier, 2) a sub-optimal (decision-fusion) classifier that treats all the measurements as independent, and 3) a sub-optimal (data-averaging) classifier that treats all the measurements as perfectly correlated. Our results indicate that the decision-fusion classifier is remarkably robust to the true statistical correlation between measurements. Thus, the decision-fusion classifier, that is the most attractive choice in view of the computational and communication constraints, is also a robust choice from a decision-theoretic viewpoint.

2 CSP Classifiers for Multiple Measurements

We consider Gaussian classifiers which assume that the underlying data has complex circular Gaussian statistics. The notation $\mathbf{x} \sim \mathcal{CN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ means that $E[\mathbf{x}] = \boldsymbol{\mu}$ and $E[\mathbf{x}\mathbf{x}^H] = \boldsymbol{\Sigma}$ and $E[\mathbf{x}\mathbf{x}^T] = \mathbf{0}$ (circular assumption). We first discuss the classifier structure for a single measurement and then generalize it to multiple measurements.

2.1 Single Measurement Classifier

Consider M target classes. Let \mathbf{x} denote a complex-valued N -dimensional feature vector corresponding to a detected event. Under hypothesis $j = 1, \dots, M$ (corresponding to j -th target class), \mathbf{x} is modeled as

$$H_j : \mathbf{x} = \mathbf{s} + \mathbf{n}, \quad j = 1, \dots, M, \quad (1)$$

where $\mathbf{s} \sim \mathcal{CN}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ denotes the Gaussian signal component corresponding to the j -th class, and $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ denotes additive white Gaussian noise. A classifier C maps the event feature vector \mathbf{x} to one of the target classes. We assume that all classes are equally likely. Thus, the optimal classifier is the maximum-likelihood (ML) classifier which takes the form [4]

$$C(\mathbf{x}) = \arg \max_{j \in \{1, \dots, M\}} p_j(\mathbf{x}) \quad (2)$$

where $p_j(\mathbf{x})$ denotes the likelihood function for j -th class which takes the following form under the complex Gaussian assumption

$$p_j(\mathbf{x}) = \frac{1}{\pi^N |\boldsymbol{\Sigma}_j + \mathbf{I}|} e^{-(\mathbf{x} - \boldsymbol{\mu}_j)^H (\boldsymbol{\Sigma}_j + \mathbf{I})^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)}. \quad (3)$$

In this paper, we assume zero-mean signals so that $\boldsymbol{\mu}_j = \mathbf{0}$ for all j and, thus, all information about the targets is contained in the covariance matrices $\boldsymbol{\Sigma}_j$. In practice, $\boldsymbol{\Sigma}_j$ has to be estimated from available training data. We assume that $\text{tr}(\boldsymbol{\Sigma}_j)$ (signal energy) is the same for all j .

2.2 Multiple Measurement Classifier

Suppose that we have K measurements (in a given modality), $\{\mathbf{x}_1, \dots, \mathbf{x}_K\}$, from different nodes available to us. We are interested in combining these measurements to achieve improved classification performance. Consider the concatenated NK -dimensional feature vector

$$\mathbf{x}^{cT} = [\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_K^T] \quad (4)$$

which has the same form as (1) under different hypotheses except for the larger number of dimensions. The noise is still white but the signal correlation matrix under H_j can be partitioned as

$$\boldsymbol{\Sigma}_j^c = \begin{bmatrix} \boldsymbol{\Sigma}_{j,11} & \boldsymbol{\Sigma}_{j,12} & \cdots & \boldsymbol{\Sigma}_{j,1K} \\ \boldsymbol{\Sigma}_{j,21} & \boldsymbol{\Sigma}_{j,22} & \cdots & \boldsymbol{\Sigma}_{j,2K} \\ \vdots & \ddots & \vdots & \vdots \\ \boldsymbol{\Sigma}_{j,K1} & \boldsymbol{\Sigma}_{j,K2} & \cdots & \boldsymbol{\Sigma}_{j,KK} \end{bmatrix} \quad (5)$$

where $\boldsymbol{\Sigma}_{j,kk'} = \text{E}[\mathbf{x}_k \mathbf{x}_{k'}^H]$ denotes the cross-covariance between the k -th and k' -th measurements. The optimal classifier operates on \mathbf{x}^c and takes the form (2) with $p_j(\mathbf{x}^c)$ given by (3) by replacing \mathbf{x} with \mathbf{x}^c and $\boldsymbol{\Sigma}_j$ with $\boldsymbol{\Sigma}_j^c$.

2.3 A Simple Measurement Model

We now present a model for measurements that is used throughout the paper. Let $K = Gn_G$. Suppose that the signal component of \mathbf{x}^c can be partitioned into G groups of n_G measurements each as

$$\mathbf{s}^{cT} = [\mathbf{s}_1^T, \dots, \mathbf{s}_1^T, \mathbf{s}_2^T, \dots, \mathbf{s}_2^T, \dots, \mathbf{s}_G^T, \dots, \mathbf{s}_G^T] \quad (6)$$

where the signal component of the n_G measurements in each group is identical and it varies independently from group to group. That is, $\{\mathbf{s}_1, \dots, \mathbf{s}_G\}$ are i.i.d. according to $\mathcal{CN}(\mathbf{0}, \mathbf{\Sigma}_j)$ under H_j . The noise measurements, on the other hand, are independent across all measurements. The above signal model can capture a range of correlation between measurements. For $K = G$ ($n_G = 1$), all the measurements have independent signal components (no correlation), whereas for $K = n_G$ ($G = 1$), all the measurements have identical signal components (maximum correlation). We consider three classifiers based on the above model.

Optimum Classifier There are two sources of classification error: background noise and the inherent statistical variability in the signals captured by $\mathbf{\Sigma}_j$'s. The optimal classifier performs signal averaging within each group to reduce the noise variance and statistical averaging over the groups to reduce the inherent signal variations. The optimum classifier operates on the NG dimensional vector

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_G \end{bmatrix} = \begin{bmatrix} \mathbf{s}_1 \\ \vdots \\ \mathbf{s}_G \end{bmatrix} + \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_G \end{bmatrix} = \mathbf{s} + \mathbf{w} \quad (7)$$

where \mathbf{y}_i are obtained by averaging the measurements in each group

$$\mathbf{y}_i = \frac{1}{n_G} \sum_{j=1}^{n_G} \mathbf{x}_{(i-1)G+j} = \mathbf{s}_i + \mathbf{w}_i, i = 1, \dots, G. \quad (8)$$

Note that \mathbf{w}_i are i.i.d. $\mathcal{CN}(\mathbf{0}, \mathbf{I}/n_G)$ due to signal averaging and \mathbf{s}_i are i.i.d. $\mathcal{CN}(\mathbf{0}, \mathbf{\Sigma}_j)$ under H_j . It can be shown that the optimal classifier takes the form

$$C_{opt}(\mathbf{y}_1, \dots, \mathbf{y}_G) = \arg \min_{j=1, \dots, M} l_{opt,j}(\mathbf{y}_1, \dots, \mathbf{y}_G) \quad (9)$$

where the (negative) log-likelihood function $l_{opt,j}(\mathbf{y})$ is given by

$$\begin{aligned} l_{opt,j}(\mathbf{y}) &= \log |\mathbf{\Sigma}_j + \mathbf{I}/n_G| + \frac{1}{G} \sum_{i=1}^G \mathbf{y}_i^H (\mathbf{\Sigma}_j + \mathbf{I}/n_G)^{-1} \mathbf{y}_i \\ &= \log |\mathbf{\Sigma}_j + \mathbf{I}/n_G| + \text{tr}((\mathbf{\Sigma}_j + \mathbf{I}/n_G)^{-1} \hat{\mathbf{\Sigma}}_y) \end{aligned} \quad (10)$$

and $\hat{\mathbf{\Sigma}}_y = \frac{1}{G} \sum_{i=1}^G \mathbf{y}_i \mathbf{y}_i^H$ is the estimated data correlation matrix of $\{\mathbf{y}_i\}$.

It is insightful to consider two limiting cases. First, suppose that $K = n_G$ ($G = 1$ - perfectly correlated measurements). In the limit of large K

$$\lim_{K \rightarrow \infty} l_{opt,j}(\mathbf{y}) = \log |\mathbf{\Sigma}_j| + \mathbf{y}_1^H \mathbf{\Sigma}_j^{-1} \mathbf{y}_1 \quad (11)$$

which shows that noise is completely eliminated and the only remaining source of error is the inherent statistical variation in the signal. Now, suppose that $K = G$ ($n_G = 1$ - independent measurements). In the limit of large K

$$\lim_{K \rightarrow \infty} l_{opt,j}(\mathbf{y}) = \log |\mathbf{\Sigma}_j + \mathbf{I}| + \text{tr}((\mathbf{\Sigma}_j + \mathbf{I})^{-1} \hat{\mathbf{\Sigma}}_y) \quad (12)$$

where $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_m + \mathbf{I}$ under H_m . In this case, all statistical variation in the signal is removed due to ensemble averaging. However, there is a bias in the estimated data correlation (relative to $\boldsymbol{\Sigma}_j$) due to noise. Both data averaging (correlated measurements) and ensemble averaging (uncorrelated measurements) contribute to improved classifier performance. However, as we will see, ensemble averaging is more critical in the case of stochastic signals.

Decision-Fusion Classifier The sub-optimal decision-fusion classifier treats all measurements as independent:

$$C_{df}(\mathbf{x}_1, \dots, \mathbf{x}_K) = \arg \min_{j=1, \dots, M} l_{df,j}(\mathbf{x}_1, \dots, \mathbf{x}_K)$$

$$l_{df,j}(\mathbf{x}) = \log |\boldsymbol{\Sigma}_j + \mathbf{I}| + \frac{1}{K} \sum_{i=1}^K \mathbf{x}_i^H (\boldsymbol{\Sigma}_j + \mathbf{I})^{-1} \mathbf{x}_i$$

$$= \log |\boldsymbol{\Sigma}_j + \mathbf{I}| + \text{tr}((\boldsymbol{\Sigma}_j + \mathbf{I})^{-1} \hat{\boldsymbol{\Sigma}}_x) \quad (13)$$

where $\hat{\boldsymbol{\Sigma}}_x = \frac{1}{K} \sum_{i=1}^K \mathbf{x}_i \mathbf{x}_i^H$ is the estimated data correlation matrix of $\{\mathbf{x}_i\}$. Note that C_{opt} and C_{df} are identical for $K = G$ in the measurement model. Note also from (13) that the M scalars $\{\mathbf{x}_i^H (\boldsymbol{\Sigma}_j + \mathbf{I})^{-1} \mathbf{x}_i\}$ for $j = 1, \dots, M$ need to be transmitted from the K nodes to the manager node. Thus, C_{df} imposes a much smaller communication (and computational) burden on the network since $M \ll N$ in general. We consider only soft decision fusion in this paper. Several other forms, including hard decision fusion, are also possible [5].

Data-Averaging Classifier The data-averaging classifier treats all measurements as correlated. It operates on the average of all measurements

$$\mathbf{y}_{da} = \frac{1}{K} \sum_{i=1}^K \mathbf{x}_i = \frac{1}{G} \sum_{i=1}^G \mathbf{y}_i = \mathbf{s}_{da} + \mathbf{w}_{da} \quad (14)$$

where $\mathbf{s}_{da} \sim \mathcal{CN}(\mathbf{0}, \boldsymbol{\Sigma}_j/G)$ under H_j and $\mathbf{w}_{da} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}/K)$ in the measurement model. The data-averaging classifier takes the form

$$C_{da}(\mathbf{y}_{da}) = \arg \min_{j=1, \dots, M} l_{da,j}(\mathbf{y}_{da})$$

$$l_{da,j}(\mathbf{y}_{da}) = \log |\boldsymbol{\Sigma}_j + \mathbf{I}/K| + \mathbf{y}_{da}^H (\boldsymbol{\Sigma}_j + \mathbf{I}/K)^{-1} \mathbf{y}_{da}. \quad (15)$$

Note that C_{opt} and C_{da} are identical for $K = n_G$. All K measurements $\{\mathbf{x}_i\}$ have to be communicated to the manager node for the computation of C_{opt} and C_{da} . However, the computational burden of C_{da} is lower than that of C_{opt} .

3 Performance Analysis of the Three Classifiers

We analyze the performance of the three classifiers for $M = 2$ classes. The analysis for $M > 2$ is more involved and is beyond the scope of this paper. Simple

union bounds can be obtained for $M > 2$ via the $M = 2$ analysis presented here. We also analyze the asymptotic performance in the limit of large number of independent measurements and also provide an entropy comparison between data and decision fusion.

We assess the performance in terms of the average probability of (correct) detection (PD)

$$PD_j = P(l_j < l_m, \forall m \neq j | H_j), \quad PD = \frac{1}{M} \sum_{j=1}^M PD_j \quad (16)$$

and the average probability of false alarm (PFA)

$$PFA_j = \frac{1}{M-1} \sum_{k=1, k \neq j}^M P(l_j < l_m, \forall m \neq j | H_k), \quad PFA = \frac{1}{M} \sum_{j=1}^M PFA_j. \quad (17)$$

For $M = 2$ the above expressions simplify to

$$\begin{aligned} PD_1 &= P(l_1 < l_2 | H_1) = 1 - PFA_2, \quad PD_2 = P(l_2 < l_1 | H_2) = 1 - PFA_1 \\ PD &= 1 - PFA. \end{aligned} \quad (18)$$

Our analysis is based on a signal model in which the covariance matrices of different targets are simultaneously diagonalizable. This model is motivated in the next section and it also simplifies the exposition to gain insight into the performance of the classifiers.

3.1 Simultaneously Diagonalizable Classes

We assume that all the covariance matrices share the same eigenfunctions

$$\boldsymbol{\Sigma}_j = \mathbf{U} \mathbf{A}_j \mathbf{U}^H, \quad j = 1, \dots, M \quad (20)$$

where \mathbf{U} represents the matrix of common (orthonormal) eigenvectors for all the classes – the different classes are characterized by the diagonal matrix of eigenvalues $\mathbf{A}_j = \text{diag}(\lambda_j[1], \dots, \lambda_j[N])$. One scenario in which this assumption is approximately valid is when the source signals for different targets can be modeled as stationary processes over the duration of the detected event. In such a case, choosing \mathbf{U} as a discrete Fourier transform (DFT) matrix would serve as an approximate set of eigenfunctions [6]. The eigenvalues will then correspond to samples of the associated power spectral densities (PSD's). The numerical results in Section 4 are based on this assumption and rely on experimental data collected in the SensIT program. Note that given the knowledge of \mathbf{A}_j in the measurement model of Section 2.3, a realization for the signal in the i -th group, from the j -th class, can be generated as

$$\mathbf{s}_i = \mathbf{U} \mathbf{A}_j^{1/2} \mathbf{z}_i, \quad \mathbf{z}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}), \quad i = 1, \dots, G. \quad (21)$$

The same \mathbf{z}_i realization is used in the i -th group and it changes independently from group to group. We assume the above signal model and analyze the classifiers in the eigen (Fourier) domain so that $\{\boldsymbol{\Sigma}_j\}$ are replaced with $\{\mathbf{A}_j\}$.

3.2 Optimal Classifier

The test statistic for the optimal classifier takes the form

$$l_{opt,j}(\mathbf{y}_1, \dots, \mathbf{y}_G) = \log |\tilde{\mathbf{A}}_j| + \frac{1}{G} \sum_{i=1}^G \mathbf{y}_i^H \tilde{\mathbf{A}}_j^{-1} \mathbf{y}_i, \quad \tilde{\mathbf{A}}_j = \mathbf{A}_j + \mathbf{I}/n_G \quad (22)$$

where \mathbf{y}_i are i.i.d. according to $\mathcal{CN}(\mathbf{0}, \tilde{\mathbf{A}}_j)$. Thus, \mathbf{y}_i can be represented as $\mathbf{y}_i = \tilde{\mathbf{A}}_j^{1/2} \mathbf{z}_i$ where $\{\mathbf{z}_i\}$ are i.i.d. $\mathcal{CN}(\mathbf{0}, \mathbf{I})$. Consider the computation of PD_1 first. It can be readily shown that under H_1

$$l_{opt,1} = \log |\tilde{\mathbf{A}}_1| + \frac{1}{G} \sum_{i=1}^G \|\mathbf{z}_i\|^2, \quad l_{opt,2} = \log |\tilde{\mathbf{A}}_2| + \frac{1}{G} \sum_{i=1}^G \mathbf{z}_i^H \tilde{\mathbf{A}}_1 \tilde{\mathbf{A}}_2^{-1} \mathbf{z}_i. \quad (23)$$

Thus,

$$PD_1 = P \left(\frac{1}{G} \sum_{i=1}^G \mathbf{z}_i^H \left[\mathbf{I} - \tilde{\mathbf{A}}_1 \tilde{\mathbf{A}}_2^{-1} \right] \mathbf{z}_i < \log |\tilde{\mathbf{A}}_2| - \log |\tilde{\mathbf{A}}_1| \right) \quad (24)$$

where the quadratic form

$$\frac{1}{G} \sum_{i=1}^G \mathbf{z}_i^H \left[\mathbf{I} - \tilde{\mathbf{A}}_1 \tilde{\mathbf{A}}_2^{-1} \right] \mathbf{z}_i = \frac{1}{G} \sum_{i=1}^G \sum_{n=1}^N |z_i[n]|^2 \left(\frac{\lambda_2[n] - \lambda_1[n]}{\lambda_2[n] + 1/n_G} \right) \quad (25)$$

is a weighted sum of $NG \chi_2^2$ random variables ($\{|z_i[n]|^2\}$) whose density and distribution functions can be analytically computed but are tedious [7]. Similarly, under H_2

$$PD_2 = P \left(\frac{1}{G} \sum_{i=1}^G \mathbf{z}_i^H \left[\mathbf{I} - \tilde{\mathbf{A}}_2 \tilde{\mathbf{A}}_1^{-1} \right] \mathbf{z}_i < \log |\tilde{\mathbf{A}}_1| - \log |\tilde{\mathbf{A}}_2| \right). \quad (26)$$

We note that the PD can be computed in closed form for $M = 2$, as indicated above, but we do not provide the explicit expression since it is rather tedious.

3.3 Decision-Fusion Classifier

The test statistic for the decision-fusion classifier takes the form

$$l_{df,j}(\mathbf{x}_1, \dots, \mathbf{x}_K) = \log |\hat{\mathbf{A}}_j| + \frac{1}{K} \sum_{i=1}^K \mathbf{x}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{x}_i, \quad \hat{\mathbf{A}}_j = \mathbf{A}_j + \mathbf{I}. \quad (27)$$

The quadratic form in the test statistic can be expanded as

$$\begin{aligned}
 \frac{1}{K} \sum_{i=1}^K \mathbf{x}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{x}_i &= \frac{1}{G n_G} \sum_{i=1}^G \sum_{k=1}^{n_G} (\mathbf{s}_i + \mathbf{n}_{(i-1)G+k})^H \hat{\mathbf{A}}_j^{-1} (\mathbf{s}_i + \mathbf{n}_{(i-1)G+k}) \\
 &= \frac{1}{G} \sum_{i=1}^G \left[\mathbf{s}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{s}_i + 2 \operatorname{Re} \left[\mathbf{s}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{w}_i \right] \right. \\
 &\quad \left. + \frac{1}{n_G} \sum_{k=1}^{n_G} \mathbf{n}_{(i-1)G+k}^H \hat{\mathbf{A}}_j^{-1} \mathbf{n}_{(i-1)G+k} \right] \quad (28)
 \end{aligned}$$

where \mathbf{s}_i and \mathbf{w}_i are defined in (8) and (21). The density for the above quadratic form can be computed exactly but here we provide a simple approximation that yields fairly accurate (but conservative) PD estimates and relates $l_{df,j}$ to $l_{opt,j}$. We make two approximations. First, we replace the $\mathbf{w}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}/n_G)$ in (28) with $\mathbf{w}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. Second, we replace $\frac{1}{n_G} \sum_{k=1}^{n_G} \mathbf{n}_{(i-1)G+k}^H \hat{\mathbf{A}}_j^{-1} \mathbf{n}_{(i-1)G+k}$ with $\mathbf{w}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{w}_i$, $\mathbf{w}_i \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. With the above approximations we have

$$\frac{1}{K} \sum_{i=1}^K \mathbf{x}_i^H \hat{\mathbf{A}}_j^{-1} \mathbf{x}_i \approx \frac{1}{G} \sum_{i=1}^G \hat{\mathbf{y}}_i^H \hat{\mathbf{A}}_j^{-1} \hat{\mathbf{y}}_i \quad (29)$$

where $\hat{\mathbf{y}}_i$ are i.i.d. $\mathcal{CN}(\mathbf{0}, \hat{\mathbf{A}}_j)$ under H_j . Thus, the PD_1 and PD_2 for the decision fusion classifier can be approximated by those of the optimal classifier given in (24) and (26) by replacing $\tilde{\mathbf{A}}_j$ with $\hat{\mathbf{A}}_j$. In particular, the quadratic form for PD_1 is given by

$$\frac{1}{G} \sum_{i=1}^G \mathbf{z}_i^H \left[\mathbf{I} - \hat{\mathbf{A}}_1 \hat{\mathbf{A}}_2^{-1} \right] \mathbf{z}_i = \frac{1}{G} \sum_{i=1}^G \sum_{n=1}^N |z_i[n]|^2 \left(\frac{\lambda_2[n] - \lambda_1[n]}{\lambda_2[n] + 1} \right) \quad (30)$$

which is a weighted sum of NG χ_2^2 random variables ($\{|z_i[n]|^2\}$), as for C_{opt} . However, the weights are different and essentially amount to a loss in SNR by a factor of n_G compared to C_{opt} since C_{df} does not do signal averaging within each group. The above conservative analysis shows that C_{df} fully exploits the independent observations across different groups, as C_{opt} , but incurs an effective loss in SNR compared to C_{opt} .

3.4 Data-Averaging Classifier

The test statistic for the data-averaging classifier takes the form

$$l_{da,j}(\mathbf{y}_{da}) = \log |\check{\mathbf{A}}_j| + \mathbf{y}_{da}^H \check{\mathbf{A}}_j^{-1} \mathbf{y}_{da} \quad , \quad \check{\mathbf{A}}_j = \mathbf{A}_j + \mathbf{I}/K \quad (31)$$

where $\mathbf{y}_{da} \sim \mathcal{CN}(\mathbf{0}, \check{\mathbf{A}}_j)$, $\check{\mathbf{A}}_j = \mathbf{A}_j/G + \mathbf{I}/K$ under H_j . Thus, \mathbf{y}_{da} can be represented as $\mathbf{y}_{da} = \check{\mathbf{A}}_j^{1/2} \mathbf{z}$ where $\mathbf{z} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. Proceeding similarly as above, it

can be shown that

$$PD_1 = P \left(\mathbf{z}^H \check{\mathbf{A}}_1 \left[\check{\mathbf{A}}_1^{-1} - \check{\mathbf{A}}_2^{-1} \right] \mathbf{z} < \log |\check{\mathbf{A}}_2| - \log |\check{\mathbf{A}}_1| \right) \quad (32)$$

where the quadratic form can be expressed as

$$\mathbf{z}^H \check{\mathbf{A}}_1 \left[\check{\mathbf{A}}_1^{-1} - \check{\mathbf{A}}_2^{-1} \right] \mathbf{z} = \sum_{n=1}^N |z[n]|^2 \left(\frac{\lambda_1[n]}{G} + \frac{1}{K} \right) \left(\frac{\lambda_2[n] - \lambda_1[n]}{(\lambda_1[n] + 1/K)(\lambda_2[n] + 1/K)} \right) \quad (33)$$

which is a weighted sum of N χ_2^2 random variables ($\{|z[n]|^2\}$). Similarly,

$$PD_1 = P \left(\mathbf{z}^H \check{\mathbf{A}}_2 \left[\check{\mathbf{A}}_2^{-1} - \check{\mathbf{A}}_1^{-1} \right] \mathbf{z} < \log |\check{\mathbf{A}}_1| - \log |\check{\mathbf{A}}_2| \right). \quad (34)$$

The density and distribution function of the quadratic form in (33) can be computed in closed form [7] and thus the PD of the data-averaging classifier can also be computed in closed form.

The data-averaging classifier provides maximum immunity against noise by averaging over all measurements. However, it does not exploit the independent signal component in different groups to reduce the inherent variations in the signal. Thus, in the limit of large number of uncorrelated measurements, we expect both C_{opt} and C_{df} to exhibit improved performance (perfect classification under certain conditions), but the performance of C_{da} will always be limited.

3.5 Asymptotic Performance

We now analyze classifier performance in the limit of large G (and K) for fixed n_G . In this analysis, we consider arbitrary $M > 2$. According to the analysis above, the only effect of n_G is to alter the effective SNR in the case of C_{opt} and C_{df} . First, consider the optimal classifier. Note that

$$l_{opt,j}(\mathbf{y}_1, \dots, \mathbf{y}_G) = -\log p_j(\mathbf{y}_1, \dots, \mathbf{y}_G)/G = -\frac{1}{G} \sum_{i=1}^G \log p_j(\mathbf{y}_i) \quad (35)$$

since \mathbf{y}_i are i.i.d. $\mathcal{CN}(\mathbf{0}, \tilde{\mathbf{A}}_m)$ under H_m . Thus, under H_m , it is well-known that by the law of large numbers [8]

$$\lim_{G \rightarrow \infty} l_{opt,j}(\mathbf{y}_1, \dots, \mathbf{y}_G) = -E_m[\log p_j(\mathbf{Y})] = D(p_m \| p_j) + h_m(\mathbf{Y}) \quad (36)$$

where $E_m[\cdot]$ denotes expectation under H_m , $D(p_m \| p_j)$ is the Kullback-Leibler distance between p_j and p_m [8]

$$D(p_m \| p_j) = E_m[\log(p_m(\mathbf{Y})/p_j(\mathbf{Y}))] = \log \left(|\tilde{\mathbf{A}}_j| / |\tilde{\mathbf{A}}_m| \right) + \text{tr} \left(\tilde{\mathbf{A}}_j^{-1} \tilde{\mathbf{A}}_m - \mathbf{I} \right) \quad (37)$$

and $h_m(\mathbf{Y})$ is the differential entropy (in bits) of \mathbf{y}_i under H_m [8]

$$h_m(\mathbf{Y}) = -\mathbb{E}_m[\log p_m(\mathbf{Y})] = \log \left((\pi\epsilon)^N |\tilde{\mathbf{A}}_m| \right). \quad (38)$$

From (36), we note that under H_m the different test statistics (for $j = 1, \dots, M$) differ only in the term $D(p_m || p_j) \geq 0$ which is identically zero for $j = m$. Thus, perfect classification ($PD = 1, PFA = 0$) is attained in the limit of large G if

$$D(p_m || p_j) > 0 \quad \forall j, m, j \neq m \quad (39)$$

which would be true in general for any given SNR (and any fixed n_G).

Now consider the decision-fusion classifier. Recall from (27) and (29) that the test statistics can be conservatively approximated as

$$l_{df,j}(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_G) \approx -\log \hat{p}_j(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_G)/G = \frac{1}{G} \sum_{i=1}^G -\log \hat{p}_j(\hat{\mathbf{y}}_i) \quad (40)$$

since $\hat{\mathbf{y}}_i$ are i.i.d. $\mathcal{CN}(\mathbf{0}, \hat{\mathbf{A}}_m)$ under H_m and \hat{p}_j denotes the density of $\mathcal{CN}(\mathbf{0}, \hat{\mathbf{A}}_j)$. Thus, in the limit of large G (under H_m)

$$\lim_{G \rightarrow \infty} l_{df,j}(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_G) = -\mathbb{E}_m[\log \hat{p}_j(\hat{\mathbf{Y}})] = D(\hat{p}_m || \hat{p}_j) + h_m(\hat{\mathbf{Y}}) \quad (41)$$

where $D(\hat{p}_m || \hat{p}_j)$ and $h_m(\hat{\mathbf{Y}})$ are defined similar to (37) and (38). Consequently, in the limit of large G we expect perfect classification if

$$D(\hat{p}_m || \hat{p}_j) > 0 \quad \forall j, m, j \neq m \quad (42)$$

which would also be true in general for any given SNR (and any fixed n_G).

Finally, consider the data-averaging classifier whose test statistics are given in (31) where $\mathbf{y}_{da} \sim \mathcal{CN}(\mathbf{0}, \check{\mathbf{A}}_m)$ under H_m . Recall that $\check{\mathbf{A}}_j = \mathbf{A}_j + \mathbf{I}/K$ and $\check{\check{\mathbf{A}}}_j = \mathbf{A}_j/G + \mathbf{I}/K$. As $G(K) \rightarrow \infty$, $\check{\check{\mathbf{A}}}_j \rightarrow \check{\mathbf{A}}_j$ and $\check{\check{\mathbf{A}}}_j \rightarrow \mathbf{0}$. Consequently,

$$\lim_{G \rightarrow \infty} l_{da,j}(\mathbf{y}_{da}) = \log |\mathbf{A}_j| \quad (43)$$

independent of the true underlying hypothesis. Thus, in the limit of large G (K), the data-averaging classifier assigns every event to the class with the smallest value of $\log |\mathbf{A}_j|$ and results in worst performance ($PD = PFA = 1/M$).

3.6 Entropy Comparison Between Data and Decision Fusion

The above analysis indicates that C_{df} approximates the performance of C_{opt} except for an SNR loss depending on the fraction of correlated measurements n_G . The numerical results in the next section confirm the analysis. However, the attractiveness of C_{df} is also implicitly based on the assumption that communicating the likelihoods from the K nodes to the manager node puts a smaller communication burden on the network compared to communicating the N -dimensional feature vectors in the case of C_{opt} .

Recall from (13) that in C_{df} the M quadratic forms $\{\mathbf{x}_i^H(\boldsymbol{\Sigma}_j + \mathbf{I})^{-1}\mathbf{x}_i, j = 1, \dots, M\}$ are communicated from the i -th node to the manager node for $i = 1, \dots, K$. In C_{opt} , on the other hand, the N -dimensional vectors \mathbf{x}_i are communicated from the K nodes to the manager node. Thus, we need to compare the cost of communicating M quadratic forms (scalars) to that of communicating an N -dimensional Gaussian vector from each node to the manager node. We compare the communication cost in terms of differential entropy [8].¹

The differential entropy of $\mathbf{x} \sim \mathcal{CN}(\mathbf{0}, \boldsymbol{\Sigma}_m + \mathbf{I})$ is

$$h_m(\mathbf{X}) = -\mathbb{E}[\log p_m(\mathbf{X})] = \log((\pi e)^N |\boldsymbol{\Sigma}_m + \mathbf{I}|) = \log((\pi e)^N |\mathbf{A}_m + \mathbf{I}|) \quad (44)$$

and quantifies the information content of any \mathbf{x}_i from m -th class.

Now consider the differential entropy of the quadratic forms used by C_{df} . Let q_{jm} denote the quadratic form associated with $l_{df,j}$ under H_m

$$q_{jm} = \mathbf{x}^H(\boldsymbol{\Sigma}_j + \mathbf{I})^{-1}\mathbf{x} = \mathbf{z}^H \mathbf{A}_{jm} \mathbf{z} \quad (45)$$

where $\mathbf{x} \sim \mathcal{CN}(\mathbf{0}, \boldsymbol{\Sigma}_m + \mathbf{I})$ and the second equality is based on the eigen-decomposition

$$(\boldsymbol{\Sigma}_m + \mathbf{I})^{1/2}(\boldsymbol{\Sigma}_j + \mathbf{I})^{-1}(\boldsymbol{\Sigma}_m + \mathbf{I})^{1/2} = \mathbf{U} \mathbf{A}_{jm} \mathbf{U}^H \quad (46)$$

which uses the representation $\mathbf{x} = (\boldsymbol{\Sigma}_m + \mathbf{I})^{1/2}\mathbf{z}$, $\mathbf{z} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$. Note that under the simultaneously diagonalizable signal model we have

$$\mathbf{A}_{jm} = \widehat{\mathbf{A}}_m \widehat{\mathbf{A}}_j^{-1}, \quad \widehat{\mathbf{A}}_j = \mathbf{A}_j + \mathbf{I}. \quad (47)$$

We want to compute the entropy of the quadratic form random variable Q_{jm} for all j, m . We first compute the worst case (highest) entropy by assuming that Q_{jm} is Gaussian. Using the fact that \mathbf{x} is Gaussian, it can be readily shown that

$$\mathbb{E}[Q_{jm}] = \mathbb{E}_m[\mathbf{x}^H(\boldsymbol{\Sigma}_j + \mathbf{I})^{-1}\mathbf{x}] = \text{tr}(\mathbf{A}_{jm}), \quad \text{var}(Q_{jm}) = \text{tr}(\mathbf{A}_{jm}^2). \quad (48)$$

Thus, the worst-case entropy of Q_{jm} is given by [8]

$$h(Q_{jm}) = \frac{1}{2} \log(2\pi e \text{tr}(\mathbf{A}_{jm}^2)). \quad (49)$$

Note from (47) that $\mathbf{A}_{jm} = \mathbf{I}$ for $j = m$. Thus, $h(Q_{jj})$ is the same for all j . For $j = m$, the true entropy can also be easily computed since $q_{jj} = \|\mathbf{z}\|^2$ from (45). Now, $q = \|\mathbf{z}\|^2 \sim \chi_{2N}^2$ with density given by [7]

$$p_Q(q) = \frac{1}{(N-1)!} q^{N-1} e^{-q}, \quad q \geq 0 \quad (50)$$

¹ We note differential entropy can be a bit misleading since it can be negative. However, a comparison of the difference in differential entropies is still valid – a quantity with higher entropy would require more bits to encode. A more intuitive interpretation of differential entropy is based on the fact that the entropy of an n -bit quantization of continuous random variable X is approximately $h(X) + n$ [8].

and thus

$$h(Q) = -\mathbb{E}[\log p_Q(Q)] = \log(N-1)! - N - (N-1) \int_0^\infty p_Q(q) \log(q) dq. \quad (51)$$

We note that the true entropy of q_{jm} , for $j \neq m$, can also be computed in closed-form but it is a bit more involved. Furthermore, as our numerical results indicate, $h(Q)$ is a good estimate for the entropy of q_{jm} for all j, m .

4 Simulation Results Based on Real Data

We now present numerical results based on real data collected in the SensIT program. We consider the problem of classifying a single vehicle. We consider $M = 2$ classes: Amphibious Assault Vehicle (AAV; tracked vehicle) and Dragon Wagon (DW; wheeled vehicle). We simulated $N = 25$ dimensional acoustic measurements from $K = Gn_G = 10$ nodes according to the model in Section 2.3. The eigenvalues (PSD samples) for the two vehicles were estimated from experimental data. The measurements at different nodes were generated using (21). The PD and PFA were estimated using Monte Carlo simulation over 5000 independent events. For $\mathbf{x} \sim \mathcal{CN}(\mathbf{0}, \mathbf{A} + \mathbf{I})$, $\text{SNR} = \text{tr}(\mathbf{A})/\text{tr}(\mathbf{I})$.

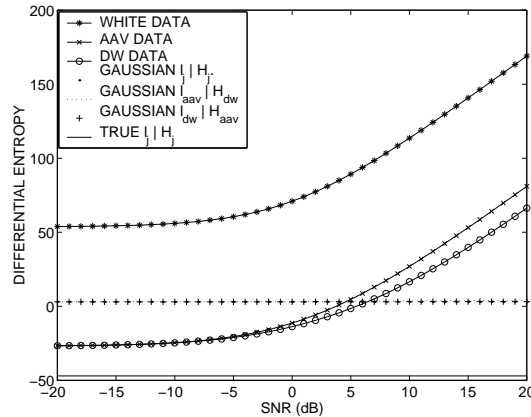


Fig. 1. Comparison between differential entropies of $N = 25$ dimensional Gaussian vectors used by l_{opt} and those of the quadratic forms used by l_{df} . The entropies for a white vector and two correlated vectors (AAV and DW) are plotted. The worst-case entropies for quadratic forms, assuming Gaussian statistics, are plotted. The true entropy of the quadratic forms, under the correct hypothesis, is also plotted. The entropy gains of decision fusion over data fusion are evident.

Figure 1 compares the differential entropy of Gaussian data in (44) with that of the quadratic forms in (49) and (51). The entropies for three data vectors are plotted: white data (maximum entropy), AAV data, and DW data. The

worst-case entropy in (49) of the quadratic forms used by C_{df} are also plotted for all j, m (they are nearly identical). It can be seen that for SNR above 5dB, the worst-case entropy of Q_{jm} is lower than that of data. The true entropy of Q_{jj} , given in (51), is also plotted for comparison. The true entropy of Q_{jj} is seen to be substantially lower compared to that of data for the entire SNR range considered. This indicates the significant potential gains of C_{df} over C_{opt} in terms of the communication burden.

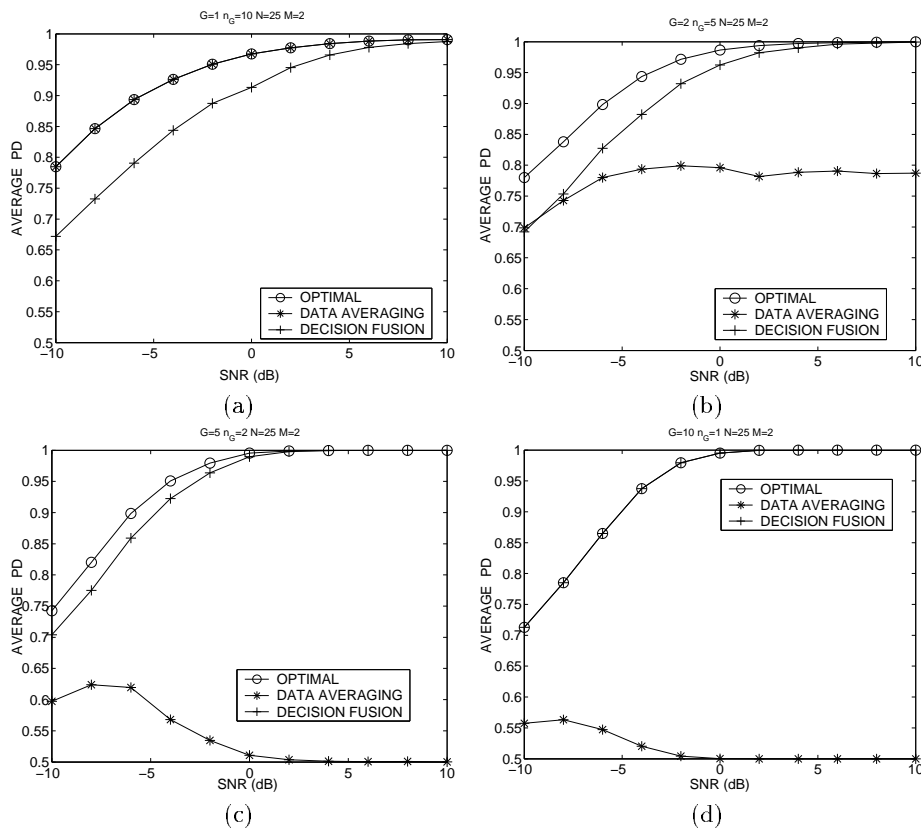


Fig. 2. PD of the three classifiers versus SNR. (a) $K = n_G = 10$ (perfectly correlated measurements). (b) $G = 2$ and $n_G = 5$. (c) $G = 5$ and $n_G = 2$. (d) $K = G = 10$ (independent measurements).

Figure 2 plots the PD as a function of SNR for the three classifiers for $K = 10$ and different combinations of G and n_G . The PFA is simply given by $1 - PD$ for $M = 2$. As expected, C_{opt} and C_{da} perform identically for $K = n_G$ (perfectly correlated case; Figure 2(a)), whereas C_{opt} and C_{df} perform identically for $K = G$ (independent case; Figure 2(d)). Note that C_{df} incurs a small loss in performance in the perfectly correlated (worst) case which diminishes at high SNRs. The performance loss in C_{da} in the independent (worst) case is very sig-

nificant and does not improve with SNR. This is consistent with our analysis. At high SNR, all events are classified as DW by C_{da} since $\log |\mathbf{A}_{DW}| < \log |\mathbf{A}_{AAV}|$ due to the peakier eigenvalue distribution for DW, as evident from Figure 3(a). Figure 3(b) compares the PD of the three classifiers for an intermediate case ($G = n_G = 2$) with $K = 4$, $N = 15$ -dimensional measurements. Analytically computed PD for C_{opt} and C_{da} and the conservative approximation for PD of C_{df} are also plotted and agree well with the simulation results.

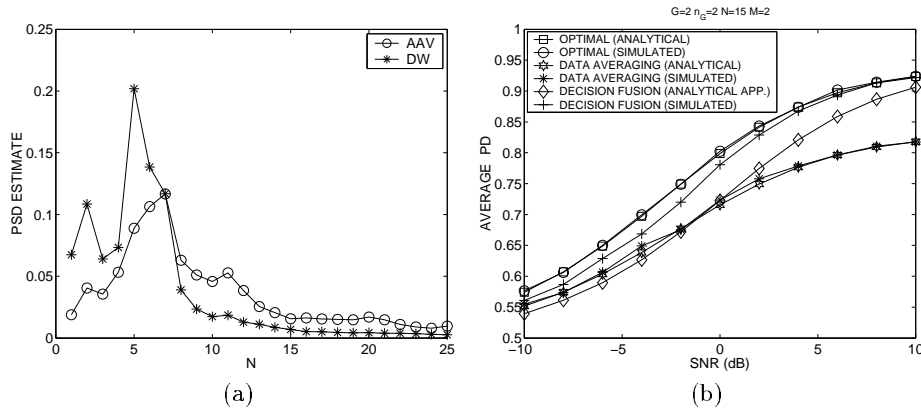


Fig. 3. (a) Covariance matrix eigenvalues (PSD estimates) for AAV and DW. (b) Comparison of simulated and analytical PD for $K = 4$, $G = 2$, $n_G = 2$ and $N = 15$.

Figure 4 plots the PD for the three classifiers as function of G ($K = 10$) for two different SNRs. It is evident that C_{df} closely approximates C_{opt} whereas C_{da} incurs a large loss when $K \neq n_G$. It is worth noting that for $\text{SNR} = -5\text{dB}$, the performance of C_{opt} and C_{df} first improves slightly with G and then gets worse again. This is consistent with the observation, in non-coherent communication over fading channels, that there is an optimal level of diversity (G) for a given SNR – increasing G beyond that level results in a loss in performance [7].

5 Conclusions

We have taken a first step in addressing the problem of how much information should be exchanged between nodes for distributed decision making in sensor networks. Our analysis is based on modeling the source signal as a stationary Gaussian process. In general, measurements from multiple nodes will provide a mixture of correlated and uncorrelated information about the source signal. The optimal classifier exploits the correlated measurements to improve the SNR and the independent measurements to stabilize the inherent statistical variability in the signal. Both effects are important for improving classifier performance. However, for stochastic signals, the fusion of independent measurements is most significant. In this context, our results demonstrate that the simple sub-optimal decision-fusion classifier, that treats all measurements as independent, is not only

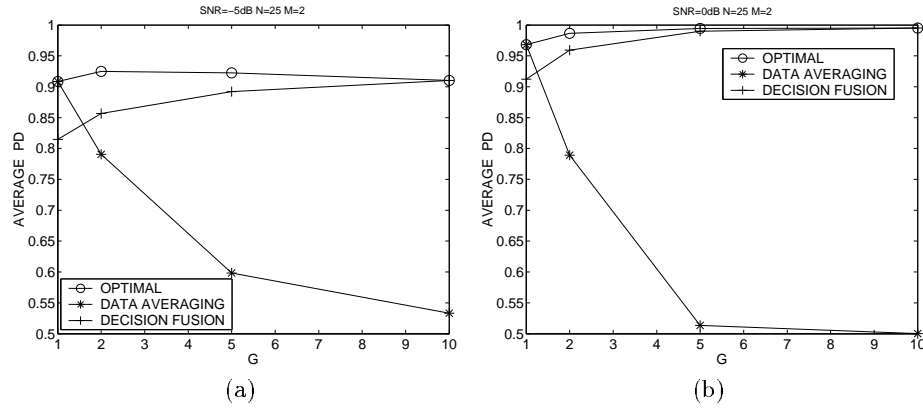


Fig. 4. Comparison of PD of the three classifiers for varying values of G ($K = 10$). (a) SNR = -5 dB. (b) SNR = 0 dB.

an attractive choice given the computational and communication constraints in a sensor network, but is also a robust choice from a decision theoretic viewpoint. The decision-fusion classifier fully exploits the independent measurements and only incurs an effective SNR loss compared to the optimal classifier depending on the fraction of correlated measurements. However, if the source signal exhibits a non-zero mean or fewer degrees of freedom (lower-rank covariance matrix), data averaging to improve SNR might become more important. We note that exploiting a non-zero mean is difficult in practice due to various sources of measurement error. Directions for future research include hard decision fusion, quantized measurements, and multiple-target classification.

References

- [1] D. Estrin, L. Girod, G. Pottie, and M. Srivastava, "Instrumenting the world with wireless sensor network," in *Proc. IEEE Int. Conf. on Acoust., Speech and Signal Proc. - ICASSP'01*, 2001, pp. 2675–2678.
- [2] "Special issue on collaborative signal and information processing in microsensor networks," in *IEEE Signal Processing Magazine*. (S. Kumar and F. Zhao and D. Shepherd (eds.)), March 2002.
- [3] D. Li, K. Wong, Y. Hu, and A. Sayeed, "Detection, classification, tracking of targets in micro-sensor networks," in *IEEE Signal Processing Magazine*, March 2002, pp. 17–29.
- [4] R. Duda, P. Hart, and D. Stork, *Pattern Classification*, Wiley, 2nd edition, 2001.
- [5] J. Kittler, M. Hatef, R. Duin, and J. Matas, "On combining classifiers," *IEEE Trans. Pattern Anal. Machine Intelligence*, vol. 20, no. 3, pp. 226–238, Mar. 1998.
- [6] Robert M. Gray, "On the asymptotic eigenvalue distribution of toeplitz matrices," *IEEE Trans. Inform. Th.*, vol. 18, no. 6, pp. 725–730, Nov. 1972.
- [7] J. G. Proakis, *Digital Communications*, McGraw Hill, New York, 3rd edition, 1995.
- [8] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, Wiley, 1991.