

# IMPROVED LOCAL DISCRIMINANT BASES USING EMPIRICAL PROBABILITY DENSITY ESTIMATION

Naoki Saito, Schlumberger-Doll Research, Ronald R. Coifman, Yale University

Naoki Saito, Schlumberger-Doll Research, Old Quarry Road, Ridgefield, CT 06877 USA

**Key Words:** Local Feature Extraction, Pattern Classification, Projection Pursuit, Density Estimation

## 1. INTRODUCTION

Recently, the authors introduced the concept of the so-called Local Discriminant Basis (LDB) for signal and image classification problems [6], [17, Chap. 4], [19], [20]. This method first decomposes available training signals in a *time-frequency dictionary* (also known as a *dictionary of orthonormal bases*) which is a large collection of the basis functions (such as wavelet packets and local trigonometric functions) localized both in time and in frequency. Then, signal energies at the basis coordinates are accumulated for each signal class separately to form a time-frequency energy distribution per class. Based on the differences among these energy distributions (measured by a certain “distance” functional), a complete orthonormal basis called LDB, which “can see” the distinguishing signal features among signal classes, is selected from the dictionary. After the basis is determined, expansion coefficients in the most important several coordinates (features) are fed into a traditional classifier such as linear discriminant analysis (LDA) or classification tree (CT). Finally, the corresponding coefficients of test signals are computed and fed to the classifier to predict their classes.

This LDB concept has been increasingly popular and applied to a variety of classification problems including geophysical acoustic waveform classification [18], radar signal classification [11], and classification of neuron firing patterns of monkeys to different stimuli [22]. Through these studies, we have found that the criterion used in the original LDB algorithm—the one based on the differences of the time-frequency energy distributions among signal classes—is not always the best one to use. Consider an artificial problem as follows. Suppose one class of signals consists of a single basis function in a time-frequency dictionary with its amplitude 10 and they are embedded in white Gaussian noise (WGN) with zero mean and unit variance. The other class of signals con-

sists of the same basis function but with its amplitude  $-10$  and again they are embedded in the same WGN process. Then their time-frequency energy distributions are identical. Therefore, we cannot select the right basis function as a discriminator. This simple counterexample suggests that we should also consider the differences of the distributions of the expansion coefficients in each basis coordinate. In this example, all coordinates except the one corresponding to the single basis function have the same Gaussian distribution. The probability density function (pdf) of the projection of input signals onto this one basis function should reveal twin peaks around  $\pm 10$ .

In this paper we propose a new LDB algorithm based on the differences among coordinate-wise pdfs as a basis selection criterion and we explain similarities and differences among the original LDB algorithm and the new LDB algorithm.

## 2. STEPS TOWARD THE ORIGINAL LDB

In this section, we review various feature extraction strategies and the original LDB algorithm. Before proceeding further, let us set up some notations. Let  $\mathcal{X} \times \mathcal{Y}$  be a set of all pairs of input signals and the corresponding class labels  $(\mathbf{x}, y)$ . We call  $\mathcal{X} \subseteq \mathbb{R}^n$  an *input signal space* and  $\mathcal{Y} = \{1, \dots, K\}$  an *output response space* which is simply a set of possible class labels (names). Most of the signal classification problems we are interested in, the dimension  $n$  is rather large. For example, in medical X-ray tomography, we typically have  $n \approx 512^2$ , and in seismic signals,  $n \approx 4000$ . A *classifier* is a function  $d : \mathcal{X} \rightarrow \mathcal{Y}$  which assigns a class name to each input signal  $\mathbf{x} \in \mathcal{X}$ . Note that if we define  $A_y \triangleq \{\mathbf{x} \in \mathcal{X} : d(\mathbf{x}) = y\}$ , then  $A_y$ s are disjoint and  $\mathcal{X} = \bigcup_{y=1}^K A_y$ . A *learning* (or *training*) dataset  $\mathcal{T}$  consists of  $N$  pairs of input signals and the corresponding class names;  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ . Let  $N_y$  be a number of class  $y$  signals in  $\mathcal{T}$ . So,  $N = \sum_{y=1}^K N_y$ . We often write  $\mathbf{x}_i^{(y)}$  to emphasize that this belongs to class  $y$ . We also assume the availability of another dataset  $\mathcal{T}'$  which is indepen-

dent of  $\mathcal{T}$  and still are sampled from the same probability model. This is called a *test* dataset and used for evaluation of classifiers. Let us now introduce a probability model. Let  $P(A, y)$  be a probability on  $\mathcal{X} \times \mathcal{Y}$  ( $A \subset \mathcal{X}$ ,  $y \in \mathcal{Y}$ ). Let  $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$  be a random sample from this probability distribution:

$$P(A, y) \triangleq P(\mathbf{X} \in A, Y = y) = \pi_y P(\mathbf{X} \in A | Y = y),$$

where  $\pi_y \triangleq P(Y = y)$  is a prior probability of class  $y$ . In practice, we often assume  $\pi_y = N_y/N$  where  $N_y$  is the number of available class  $y$  samples in the training dataset  $\mathcal{T}$  (i.e.,  $\sum_y N_y = N$ ).

The first and most direct approach to the classification problem seems:

**Approach 0:** *Construct the best possible classifier using information of  $P(A, y)$ .*

This naturally leads to the concept of the *Bayes classifier*, but we quickly realize that it is impossible to obtain in practice.

## 2.1. The Bayes classifier

If we know the true probability distribution  $P(A, y)$ , then the optimal classifier is the so-called *Bayes classifier* (or *rule*)  $d_B$  and is defined as

$$P(d_B(\mathbf{X}) \neq Y) \leq P(d(\mathbf{X}) \neq Y),$$

for any other classifier  $d(\mathbf{X})$ . Then, the *Bayes* misclassification rate  $R_B \triangleq P(d_B(\mathbf{X}) \neq Y)$  is clearly the minimum among the rates obtained by all possible classifiers using information of  $P(A, y)$ . In fact, assuming the existence of the conditional pdf  $p(\mathbf{x} | y)$  such that  $P(A, y) = \int_A p(\mathbf{x} | y) d\mathbf{x}$ , it is well known (see e.g., [3, p. 14]) that the Bayes classifier is to “assign  $\mathbf{x}$  to class  $k$  if  $\mathbf{x} \in A_k^*$ ,” where  $A_k^* = \{\mathbf{x} \in \mathcal{X} : p(\mathbf{x} | k)\pi_k = \max_{y \in \mathcal{Y}} p(\mathbf{x} | y)\pi_y\}$ . This Bayes classifier, however, is impossible to construct. First of all,  $p(\mathbf{x} | y)$  is not known in practice. It is even difficult to estimate  $p(\mathbf{x} | y)$  computationally using available training samples because of the high dimensionality of the input space  $\mathcal{X}$  (*curse of dimensionality*); we need a huge number of training samples to get reliable estimate of  $p(\mathbf{x} | y)$ . Therefore, we need to reduce the dimensionality of the problem without losing important information for classification.

## 2.2. Feature extraction, dimension reduction, and projection pursuit

Faced with the curse of dimensionality and having such difficulty in constructing the Bayes classifiers, the extraction of important features becomes essential. We want to keep only important information and discard the irrelevant information for classification purposes. Humans perform this kind of “feature compression” on a daily basis when facing classification and discrimination tasks. As Scott mentions in his book [21, Chap. 7], this strategy is also supported by the empirical observation that multivariate data in  $\mathbb{R}^n$  are almost never  $n$ -dimensional and there often exist lower dimensional structures of data. In other words, a signal classification problem often has an *intrinsic dimension*  $m < n$ . (Note that this is clearly different from that of signal compression problem even for the same dataset; the intrinsic dimension is goal dependent.) Therefore, it would be much more efficient and effective to analyze the data and build classifiers in this smaller dimensional subspace  $\mathcal{F}$  of  $\mathcal{X}$ , if possible. We call  $\mathcal{F}$  a *feature space*, and a map  $f : \mathcal{X} \rightarrow \mathcal{F}$  a *feature extractor*. Then, the key is how to construct this “good” feature space  $\mathcal{F}$  consisting of discriminant features and design the corresponding feature extractor  $f$ . If we precisely know the underlying physical and mathematical models of the problem, then we can design a mechanism to extract specific features relevant for that problem and may obtain its intrinsic dimension. It is often difficult, however, to set up exact mathematical models for the problem we are interested in, such as medical and geophysical diagnostics. Therefore, we want to adopt exploratory approaches; we want to find discriminant features by automatic procedures and examine their effectiveness. In turn, this may lead to our understanding of the underlying physics of the problem.

Based on the above observations, the next approach to signal classification problems can be stated as follows.

**Approach 1:** *Extract the best possible features from the signals and construct the best possible classifier on these features.*

This approach can be symbolically written as:

$$d = g \circ f, \tag{1}$$

where  $f$  is a feature extractor and  $g : \mathcal{F} \rightarrow \mathcal{Y}$  is a classifier using signal features. Ideally, mapping input signals to  $\mathcal{F}$  should reveal separate clusters of points corresponding to signal classes and ease the subsequent classification tasks. In general,  $f$  can be nonlinear, e.g., neural network feature extractors such as self-organizing maps

[14]. Although this type of feature extractors have interesting capabilities, they almost always require a nonlinear optimization procedure which is computationally expensive and is difficult to find the global optimum. Then, how about using some linear transforms as  $f$ ? This question naturally leads to the concept of projection pursuit.

The *projection pursuit* (PP), which was originally proposed by Kruskal [15] and was implemented, experimented, and named by Friedman and Tukey [10], is one of the few techniques to do the dimension reduction. The original purpose of PP was to pick “interesting” low-dimensional projections of high-dimensional point clouds automatically by numerically optimizing a certain objective function or *projection index*. One can find a sequence of best one-dimensional projections by optimizing the projection index, then removing the structure that makes this direction interesting, and iterating. As one can see, the idea of PP is very general and was extended for various purposes including density estimation, regression, classification, and clustering; see excellent expository papers by Huber [12] and by Jones and Sibson [13]. In particular, by changing the projection index to the appropriate ones, many of the classical multivariate data analysis techniques, such as principal component analysis (PCA) and linear discriminant analysis (LDA), are shown to be the special cases of PP. Therefore, PP with a projection index measuring discriminant power of projections (or coordinates) seems an attractive approach for feature extraction and dimension reduction. The problem, however, still exists: 1) A straightforward application of PP may still be computationally too expensive for high dimensional problems. 2) Sequentially obtaining the best 1D projections may be too “greedy.” It may miss the important 2D structures. After all, “the best two measurements are not the two best” [7]. 3) Interpretation of the results often becomes difficult because the best 1D projection squeezes all discriminant features in the signals—however separated in time and frequency domains—into a single feature vector.

### 2.3. The original LDB method

Faced with the above difficulties, one may wonder how to approach the problem. Here, the so-called *time-frequency dictionary* comes in. This is a large collection of specific basis vectors in  $\mathcal{X}$  organized in a hierarchical binary tree structure. This dictionary is redundant and contains up to  $n(1 + \log_2 n)$  basis vectors (also called *time-frequency atoms*) which have specific time-frequency localization properties. Examples include wavelet packets and local trigonometric functions

both of which were created by Coifman and Meyer; see [16], [24] and references therein for the detailed properties of these basis functions. See also [18] in this proceeding. This dictionary of bases is huge; it contains more than  $2^n$  different complete orthonormal bases [24, p.256]. This dictionary, however, offers at least two things for us: 1) Efficiency in representing features localized both in time and in frequency. This is particularly useful for representing nonstationary signals or discontinuous signals. Moreover, this property makes interpretation of classification results far easier than Approach 0 and 1. 2) Computational efficiency. Expanding a signal of length  $n$  into such a tree-structured bases is fast, i.e.,  $O(n \log n)$  for a wavelet packet dictionary and  $O(n[\log n]^2)$  for a local trigonometric dictionary. Moreover, the Coifman-Wickerhauser best-basis algorithm [24, Chap. 8] allows one to search a “good” basis for her needs in  $O(n)$ .

The best-basis algorithm of Coifman and Wickerhauser, however, was developed mainly for signal compression problems. The LDB algorithm was developed to fully utilize these properties of the time-frequency dictionaries and the best-basis algorithm for signal classification and discrimination problems [6], [17, Chap. 4], [19], [20]. It is much more specific than Approach 0 and 1, but it offers a computationally efficient dimension reduction and extraction of local signal features. It is also “modest” in the sense that it picks a set of good coordinates from a finite collection of orthonormal bases rather than a sequence of the absolutely best 1D projections without constraint. This philosophy can be phrased as

**Approach 2:** *Extract the most discriminant features from a time-frequency dictionary, construct several classifiers on these features, and pick the best one among them.*

For a specific classifier  $g$  (which can be any conventional classifier of one’s choice, such as LDA, CT, or artificial neural networks, etc.), this approach can be written as

$$d = g \circ \Theta_m \circ \Psi. \quad (2)$$

The feature extractor here consists of  $\Psi \in O(n)$ , an  $n$ -dimensional orthogonal matrix (i.e., an orthonormal basis) selected from a time-frequency dictionary, and a *feature selector*  $\Theta_m$  selects the most important  $m (< n)$  coordinates (features) from  $n$ -dimensional coordinates. Most statistical literature focuses on the performance and statistical properties of various classifiers  $g$  in (1), (2). Some literature discusses the feature selector  $\Theta_m$  given a set of features. On the other hand, both the original and the new LDB methods focus on  $f$ , in particular, how to select  $\Psi$  from a finite collection of bases.

Let  $\mathfrak{D}$  be this time-frequency dictionary consisting of a collection of basis vectors  $\{\mathbf{w}_i\}$ ,  $i = 1, \dots, n(1 + \log_2 n)$ .  $\mathfrak{D}$  can also be expressed as a list of all possible orthonormal bases  $\{B_j\}$ ,  $j = 1, \dots, M$  with  $M > 2^n$ . Let  $B_j = \{\mathbf{w}_{j_1}, \dots, \mathbf{w}_{j_n}\}$  and let  $\mathcal{D}(B_j)$  be a measure of efficacy of the basis  $B_j$  for classification tasks. Then, both the original and the new LDB are selected rapidly by the best-basis algorithm using the following criterion:

$$\Psi = \operatorname{argmax}_{B_j \in \mathfrak{D}} \mathcal{D}(B_j). \quad (3)$$

The heart of a matter is this measure of efficacy  $\mathcal{D}$ . Here we describe the one adopted in the original LDB method. In the next section, we consider the various possibilities of this measure and propose the new LDB method.

Let  $\Gamma^{(y)}(\mathbf{w}_j)$  be a normalized total energy of class  $y$  signals along the direction  $\mathbf{w}_j$ :

$$\Gamma^{(y)}(\mathbf{w}_j) \triangleq \frac{\sum_{i=1}^{N_y} |\mathbf{w}_j \cdot \mathbf{x}_i^{(y)}|^2}{\sum_{i=1}^{N_y} \|\mathbf{x}_i^{(y)}\|^2}, \quad (4)$$

where  $\cdot$  denotes the standard inner product in  $\mathbb{R}^n$ . We refer to the tree-structured set of normalized energies  $\{\Gamma^{(y)}(\mathbf{w}_j) : \mathbf{w}_j \in \mathfrak{D}\}$  as the *normalized time-frequency energy map* of class  $y$ . Let  $\mathbf{\Gamma}^{(y)}(B_j) = (\Gamma^{(y)}(\mathbf{w}_{j_1}), \dots, \Gamma^{(y)}(\mathbf{w}_{j_n}))$  be a vector of such normalized energies corresponding to the basis  $B_j$  which is extracted from the time-frequency energy map. In the original LDB algorithm, as the measure  $\mathcal{D}$ , we have examined several functionals measuring “distances” among  $K$  vectors,  $\mathbf{\Gamma}^{(1)}(B_j), \dots, \mathbf{\Gamma}^{(K)}(B_j)$ , such as relative entropy, Hellinger distances, and simple  $\ell^2$  distances. See Equations (2)–(5) in [18] where we compared their classification performance. We note that the functionals such as relative entropy and Hellinger distance were originally created to measure the “distances” or deviations among pdfs [1], and in the original LDB algorithm, the normalized time-frequency energy maps are viewed as pdfs and the distances among them are computed for basis evaluation.

### 3. THE NEW LOCAL DISCRIMINANT BASES

In this section, we use the probability model to reinterpret the original LDB algorithm and then propose the new basis selection criteria.

#### 3.1. Discriminant measures for 1D projections

Let us first reconsider what is a good 1D projection for classification and discrimination in general. The more

specific issues for the LDB methods are treated in the next subsection.

If we project an input signal  $\mathbf{X} \in \mathcal{X}$  onto a unit vector  $\mathbf{w}_i \in \mathfrak{D}$ , then its projection (or coordinate)  $Z_i \triangleq \mathbf{w}_i \cdot \mathbf{X}$  is also a random variable. We also use the notation  $Z_i^{(y)}$  if we want to emphasize the projection of class  $y$  signals. We are interested to know how  $Z_i$  is distributed for each signal class so that we can quantify the efficacy of the direction  $\mathbf{w}_i$  for classification. We refer to such a projection index as a *discriminant measure*. We also use the term *discriminant power* of  $\mathbf{w}_i$  which is the actual value of such a measure evaluated at  $\mathbf{w}_i$ . We can think of four possibilities of such a measure:

**Type I:** A measure based on the differences of derived quantities from projection  $Z_i$ , such as mean class energies or cumulants.

**Type II:** A measure based on the differences among the pdfs of  $Z_i$ .

**Type III:** A measure based on the differences among the cumulative distribution functions (cdfs) of  $Z_i$ .

**Type IV:** A measure based on the actual classification performance (e.g., a rate of correct classification) using the projection of the available training signals.

Type I measure is the one used in the original LDB method. As we will show later, this is a special case of Type II measure. Type II measure is the one we adopt for the new LDB method. This requires the estimation of the pdfs of  $Z_i$  for each class and we will describe the details below. Type III and IV measures are currently under investigation. J. Buckheit suggested using Type III measures since computing an empirical cdf is simpler and easier than estimating a pdf. This includes the Anderson-Darling distance which he used for measuring non-Gaussianity of projections in his paper [4]. As for Type IV measures, there are two approaches. One is to estimate the 1D pdfs, assume them as the “true” 1D pdfs, and invoke the Bayes classifier to get the misclassification rate. The other is to use any simple classifier (e.g., LDA, CT,  $k$ -nearest neighbor, etc.) or a combination of them to get the least misclassification rate. The Type IV strategy is essentially the same as the local regression basis proposed by the authors [6], [17, Chap. 5]. We emphasize that the misclassification rate in this context is simply used to evaluate the direction  $\mathbf{w}_i$ , and is different from the final misclassification rate obtained by the entire system  $d$ .

In the following, we focus on Type II measures and its relationship to Type I measures. If we knew the “true”

conditional pdf  $p(\mathbf{x} | y)$ , then, the pdf of  $Z_i$  for class  $y$  signals would be:

$$q(z | y, \mathbf{w}_i) \triangleq \int_{\mathbf{w}_i \cdot \mathbf{x} = z} p(\mathbf{x} | y) d\mathbf{x}. \quad (5)$$

For notational convenience, we write  $q_i^{(y)}(z)$  instead of  $q(z | y, \mathbf{w}_i)$ . In practice, however, (5) cannot be computed since  $p(\mathbf{x} | y)$  is not available as mentioned in the previous section. Therefore, we must estimate  $q_i^{(y)}(z)$  empirically using the available training dataset. Let  $\hat{q}_i^{(y)}(z)$  denote such an estimate. There are many empirical density estimation techniques as described in [21]. In fact, the pdf estimation using wavelets and their relatives are a quite interesting subject itself [8]. In Section 5, we use a particular estimation method called *averaged shifted histograms* (ASH) [21, Chap. 5] which is computationally fast and has an interesting connection with the ‘‘spin cycle’’ algorithm of Coifman and Donoho [5] for signal denoising problems.

Having estimated pdfs,  $\hat{q}_i^{(1)}, \dots, \hat{q}_i^{(K)}$ , what kind of discriminant measure  $\mathcal{D}$  should we use to evaluate the direction  $\mathbf{w}_i$ ? The worst direction is the one which makes all the pdfs look identical and for this direction the value of  $\mathcal{D}$  should be 0. The best direction is the one which makes all the pdfs look most different or ‘‘distant’’ from one another which in turn should ease the subsequent classification tasks. Therefore  $\mathcal{D}$  should attain a maximum positive value at that direction. Also, we should have  $\mathcal{D}(\hat{q}_i^{(1)}, \dots, \hat{q}_i^{(K)}) \geq 0$  and the equality should hold if and only if  $\hat{q}_i^{(1)} \equiv \dots \equiv \hat{q}_i^{(K)}$ . These conditions indicate that  $\mathcal{D}$  should be a ‘‘distance-like’’ quantity among pdfs. There are many ‘‘distance’’ measures to quantify the difference or discrepancy of pdfs; see e.g., [1]. Among them, two popular distance measures (between two pdfs  $p$  and  $q$ ) are:

- Relative entropy (also known as cross entropy and Kullback-Leibler divergence):

$$D(p, q) \triangleq \int p(x) \log \frac{p(x)}{q(x)} dx. \quad (6)$$

- Hellinger distance:

$$H(p, q) \triangleq \int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx. \quad (7)$$

We also note that relative entropy (6) is not a metric since it satisfies neither symmetry nor triangle inequality. For measuring discrepancies among  $K$  pdfs,  $p^{(1)}, \dots, p^{(K)}$ ,

the simplest approach is to take  $\binom{K}{2}$  pairwise combinations of  $\mathcal{D}$ :

$$\mathcal{D}(p^{(1)}, \dots, p^{(K)}) \triangleq \sum_{i=1}^{K-1} \sum_{j=i+1}^K \mathcal{D}(p^{(i)}, p^{(j)}). \quad (8)$$

**Remark 3.1.** For a small number of classes, say 2 to 4, this approach may be sufficient. The more signal classes one has in her problem, the more obscure the meaning of (8) becomes; a large value of (8) may be due to a few significant terms with negligible majority (a favorable case) or to the accumulation of many terms with relatively small values (an unfavorable case). For such a case, we can take a different approach suggested by Watanabe and Kaminuma [23]. Instead of constructing a single classifier for the entire problem, consider  $K$  sets of two class problems by splitting the training dataset into class  $y$  and *non-class*  $y$ . Then construct a classifier for each two class problem. Suppose the probability of a signal  $\mathbf{x}$  being classified to class  $y$  using the  $y$ th classifier is  $p^{(y)}(y | \mathbf{x})$ . Then we assign  $\mathbf{x}$  to class  $k$  where  $k = \operatorname{argmax}_{y \in \mathcal{Y}} p^{(y)}(y | \mathbf{x})$ . We are currently investigating this approach in the LDB setting and comparing with the one based on (8).

Now, let us consider a particular Type I measure, the normalized energy (or the normalized second moment) of signals along the direction  $\mathbf{w}_i$ . This quantity for class  $y$  signals can be written as

$$V_i^{(y)} \triangleq \frac{E[Z_i^2 | Y = y]}{\sum_{i=1}^n E[Z_i^2 | Y = y]}. \quad (9)$$

Therefore, once we get the estimate  $\hat{q}_i^{(y)}$ , we can estimate the energy as well:

$$\hat{V}_i^{(y)} = \frac{\int z^2 \hat{q}_i^{(y)}(z) dz}{\sum_{i=1}^n \int z^2 \hat{q}_i^{(y)}(z) dz}. \quad (10)$$

This shows that this Type I measure can be derived once we get the estimate of pdfs. This is true for other derived quantities such as cumulants. Now, how are (9) and (10) related with the normalized time-frequency energy map used in the original LDB algorithm? With the available finite samples in the training dataset, we have

$$E[Z_i^2 | Y = y] \approx \frac{1}{N_y} \sum_{j=1}^{N_y} |\mathbf{w}_i \cdot \mathbf{x}_j^{(y)}|^2 = \int z^2 \hat{q}_i^{(y)}(z) dz.$$

Also, because of the orthonormality of  $\{\mathbf{w}_i\}$ ,

$$\begin{aligned} \sum_{i=1}^n E[Z_i^2 | Y = y] &= E[\|\mathbf{X}\|^2 | Y = y] \\ &\approx \frac{1}{N_y} \sum_{j=1}^{N_y} \|\mathbf{x}_j^{(y)}\|^2. \end{aligned}$$

Thus from (4), we conclude

$$\Gamma^{(y)}(\mathbf{w}_i) = \hat{V}_i^{(y)}.$$

In other words, the normalized time-frequency energies used in the original LDB algorithm turns out to be the finite sample version of the Type I measure (9) under the probabilistic setting. For a fixed  $\mathbf{w}_i$ , both  $V_i^{(y)}$  and  $\Gamma_i^{(y)}$  are point estimates, and the same is true of the discriminant measure  $\mathcal{D}(\Gamma^{(1)}(\mathbf{w}_i), \dots, \Gamma^{(K)}(\mathbf{w}_i))$ . On the other hand, the discriminant measure  $\mathcal{D}(\hat{q}_i^{(1)}, \dots, \hat{q}_i^{(K)})$  uses the entire probability distribution information. Therefore, the Type II measures using the pdfs can capture more subtle discriminant information such as *phase* information than the Type I measures based on (9) and (4).

**Remark 3.2.** Once noticed that the quantity (9) is essentially a point estimate of the nonlinear function (square function in this case) of a random variable  $Z_i$ , we may further improve the original LDB by the use of other nonlinear transformations of the coordinate values. That is, instead of using  $Z_i^2$  in (9) we can use another nonlinear function  $\eta(Z_i)$ , e.g.,  $\arctan Z_i$ ,  $\log(1 + |Z_i|)$ , and so on. We can also use this idea for the discriminant measure  $\mathcal{D}$  by estimating the pdf of  $\eta(Z_i)$  instead of that of  $Z_i$ .

### 3.2. Discriminant power of bases in dictionaries

Instead of just a single direction, suppose we are given a basis  $B = (\mathbf{w}_1, \dots, \mathbf{w}_n)$  in a time-frequency dictionary. Because it is difficult to estimate the  $n$ -dimensional pdf directly for a large  $n$ , we evaluate each basis vector separately and sum up their discriminant powers. For notational convenience, let  $\delta_i \triangleq \mathcal{D}(\hat{q}_i^{(1)}, \dots, \hat{q}_i^{(K)})$ , i.e., discriminant power of a single direction  $\mathbf{w}_i$ . Then, a simple-minded measure of the discriminant power of  $B$  may be

$$\mathcal{D}_n(B) \triangleq \sum_{i=1}^n \delta_i$$

However, once we start comparing the discriminant power of various bases in the dictionary, we quickly realize the shortcoming of this simple measure: many small  $\delta_i$ s may add up to a large discriminant power of the basis,

as mentioned in Remark 3.1. Therefore, we want to sum only  $k (< n)$  largest terms, i.e.,

$$\mathcal{D}_k(B) \triangleq \sum_{i=1}^k \delta_{(i)}. \quad (11)$$

where  $\{\delta_{(i)}\}$  is the decreasing rearrangement of  $\{\delta_i\}$ . In fact, it may be better not to take the most useless (i.e., non-discriminant) vectors into consideration. However, the automatic choice of  $k$  is not necessarily easy and needs further research.

Another possibility is to measure only the discriminant powers of the directions that carry the signal energies larger than a certain threshold  $t > 0$ :

$$\mathcal{D}'_t(B) \triangleq \sum_{i=1}^n \varepsilon_i \delta_i, \quad (12)$$

where  $\varepsilon_i = 1$  if  $E[Z_i^2] > t$  and  $= 0$  otherwise. The selection of  $t$  should be done carefully. It should be large enough to remove all the noisy coordinates, but also should be small enough not to discard the subtle discriminant features which are not noise.

**Remark 3.3.** It is interesting to use joint distribution of pairs of projections  $(\mathbf{w}_i, \mathbf{w}_j)$  instead of 1D projections because hidden structures (such as holes) may be captured by a 2D projection but not by any 1D projection. Since it is too expensive to estimate all possible combinations, i.e.,  $\binom{n}{2}$  joint pdfs, for each class, we select the most discriminant  $k (< n)$  coordinates in  $B$  using the 1D projections, then deal with  $\binom{k}{2}$  joint pdfs. Two-dimensional joint pdfs are indispensable for feature extraction using a dictionary of local Fourier bases. There, each expansion coefficient is a complex number. Taking magnitudes of the coefficients in this case clearly ignores important *phase* information. We are currently investigating this approach.

### 3.3. The new LDB selection algorithm

We summarize the new procedure here with its computational cost.

**Step 0:** Expand each training signal  $\mathbf{x}_j$  into a specified time-frequency dictionary  $\mathcal{D}$ . [ $O(n \log_2 n)$  or  $O(n(\log_2 n)^2)$ ]

**Step 1:** For each vector  $\mathbf{w}_i$  in the dictionary, estimate the pdfs  $\hat{q}_i^{(y)}$  of the projection,  $y = 1, \dots, K$ . [ $O(n(1 + \log_2 n))$  using ASH estimates]

**Step 2:** For each vector  $w_i$  in the dictionary, compute its discriminant power  $\delta_i$ . [ $O(n \log_2 n)$ ]

**Step 3:** Evaluate each basis  $B$  in the dictionary and obtain the best basis via

$$\Psi = \underset{B \in \mathcal{D}}{\operatorname{argmax}} \mathcal{D}(B),$$

where  $\mathcal{D}$  is either  $\mathcal{D}_k$  in (11) or  $\mathcal{D}'_t$  in (12). [ $O(n)$ ]

**Step 4:** Select  $m$  vectors from  $\Psi$  corresponding to the  $m$  largest  $\delta_i$ .

**Step 5:** Construct classifiers  $g$  with the features derived from  $m$  vectors.

We note that it is not necessary to have  $k = m$  in general.

#### 4. “SPIN CYCLE”: A DATA STABILIZATION PROCEDURE FOR TIME-FREQUENCY DICTIONARIES

Our bases in the dictionaries, i.e., wavelets, wavelet packets, and local trigonometric bases, do not have the translation invariance property: if we shift the original signal, its expansion coefficients change (significantly in some cases) from those of the original signal so that one cannot tell the amount of shift from the coefficients. Depending on the applications, this lack of translation invariance may be problematic. For example, for image texture classification, we do not concern the locations of individual edges or transients which form texture elements. Thus, it is preferable to make the analysis and classification processes more insensitive to translations in such problems. On the other hand, if the time delay of a certain waveform is an important discriminant feature in one’s problem, then the lack of translation invariance may not be too critical.

For compensating the lack of translation invariance, we use the so-called *spin cycle* procedure: increase the number of sample signals in the training and the test datasets by creating their translated versions. More precisely, we shift each  $x_i$  in the training and the test datasets in a circular manner by  $-\tau, -\tau + 1, \dots, -1, 1, \dots, \tau$ , where  $\tau \in \mathbb{N}$  and  $\tau < n$ . Then we have  $2\tau + 1$  signals for each original signal (counting itself) all of which share the same class assignment  $y_i$ . Next, we construct the LDB using this increased training dataset, extract top  $m$  features, and build a classifier. Then we feed all the signals to the classifier and predict the class labels. Finally, for each original signal, we take the majority vote on the  $2\tau + 1$  predicted class labels. This “spin cycle”

procedure also plays an important role for other applications such as denoising [17, Chapter 3], [5].

**Remark 4.1.** It turns out that increasing the number of sample signals by the spin cycle procedure is, in spirit, very similar to the “bagging” (*bootstrap aggregating*) procedure proposed by Breiman [2]. This method tries to stabilize certain classifiers by: 1) generating multiple versions of training dataset by the bootstrap method [9], 2) constructing a classifier for each training dataset, and 3) predicting the class of test samples by the majority vote on the predictions by all the classifiers.

## 5. EXAMPLES

In this section, we analyze two classification problems. The first one is the famous “waveform” classification problem described in the CART book [3]. The second problem is a discrimination of geophysical acoustic waveforms propagated through different media [18].

**Example 5.1.** *Triangular waveform classification.*

This is a three-class classification problem using synthetic triangular waveform data. The dimensionality of the signal was extended from 21 in [3] to 32 for the dyadic dimensionality requirement of the bases under consideration. We generated 100 training signals and 1000 test signals for each class by the following formula:

$$\begin{aligned} x^{(1)}(i) &= uh_1(i) + (1-u)h_2(i) + \epsilon(i) && \text{for Class 1,} \\ x^{(2)}(i) &= uh_1(i) + (1-u)h_3(i) + \epsilon(i) && \text{for Class 2,} \\ x^{(3)}(i) &= uh_2(i) + (1-u)h_3(i) + \epsilon(i) && \text{for Class 3,} \end{aligned}$$

where  $i = 1, \dots, 32$ ,  $h_1(i) = \max(6 - |i - 7|, 0)$ ,  $h_2(i) = h_1(i - 8)$ ,  $h_3(i) = h_1(i - 4)$ ,  $u$  is a uniform random variable on the interval  $(0, 1)$ , and  $\epsilon(i)$  are the standard normal variates. We conducted three sets of different classification experiments.

1. The original LDB method with and without Spin Cycle.
2. The new LDB method with and without Spin Cycle using the simple histograms as the pdf estimation method.
3. The new LDB method with and without Spin Cycle using ASH as the pdf estimation method.

In each experiment, we repeated the whole process 10 times by generating 10 different realizations of the training and test datasets. As a classifier  $g$  in (2), we used linear discriminant analysis (LDA) and classification tree

Method (Coordinates)	Training	Test
LDA on STD	12.0 %	22.7 %
LDA on OLDB5	14.1 %	16.2 %
LDA on OLDB5SC3	14.1 %	16.1 %
LDA on OLDB5SC5	16.2 %	17.4 %
CT on STD	7.0 %	29.3 %
CT on OLDB5	8.1 %	21.9 %
CT on OLDB5SC3	5.8 %	21.4 %
CT on OLDB5SC5	7.7 %	22.5 %

Table 1: The average misclassification rates of the waveform classification example over 10 simulations. In Method column, STD, OLDB5, represent the standard coordinates and the top 5 coordinates of the original LDB (based on the time-frequency energy distributions), respectively. SC3, SC5 mean the Spin Cycle with 3 and 5 shifts ( $\tau = 1, 2$ )

Method (Coordinates)	Training	Test
LDA on NLDB5	17.0 %	19.6 %
LDA on NLDB5SC3	16.2 %	18.2 %
LDA on NLDB5SC5	17.5 %	18.8 %
CT on NLDB5	9.9 %	26.3 %
CT on NLDB5SC3	6.3 %	24.5 %
CT on NLDB5SC5	7.9 %	23.7 %

Table 2: The average misclassification rates with the new LDB algorithm using the simple histograms as the empirical pdf estimation method (averaged over 10 simulations). Here, NLDB5xxx means that  $m = k = 5$  in (11).

(CT). As a dictionary for LDB, we used the wavelet packet dictionary with the 6-tap coiflet filter [24, Appendix C]). For the discriminant measure, we adopted the relative entropy (8) for three classes (see also (8)). As for  $m$ , the number of most important features to be fed to classifiers, we set  $m = 5$  by heuristics,  $m \approx 0.1n$  or  $0.2n$ . For comparison, we also conducted the direct application of LDA and CT over the signals represented in the standard (or canonical) basis of  $\mathbb{R}^{32}$ .

The averaged misclassification rates are summarized in Tables 1,2,3.

We would like to note that according to Breiman et al. [3], the Bayes error of this example is about 14 %. From these tables, we observe:

- The best result so far was obtained by applying LDA to NLDB5SC3, i.e., the top 5 LDB coordinates with ASH estimates on the Spin Cycled data (with three

Method (Coordinates)	Training	Test
LDA on NLDB5	16.0 %	18.2 %
CT on NLDB5	8.8 %	24.5 %
LDA on NLDB5SC3	14.4 %	<b>15.9 %</b>
CT on NLDB5SC3	5.5 %	20.8 %
LDA on NLDB5SC5	15.9 %	17.9 %
CT on NLDB5SC5	7.5 %	22.8 %

Table 3: The average misclassification rates with the new LDB algorithm using ASH as the empirical pdf estimation method (averaged over 10 simulations).

shifts).

- The four of the top 5 original and new LDB vectors are the same except their order of importance.
- Due to the nature of the problem (three triangles are located at the fixed positions), Spin Cycle with too many shifts (five shifts here) degraded the performance.
- Over all the LDA gave much better results than CT.
- For the new LDB algorithm, the Spin Cycle is critical. Without Spin Cycle, it is worse than the original LDB.
- Over all ASH-based methods gave the better results than the simple histogram-based methods.

Note that the best results NLDB5SC3 essentially used the *double* stabilization (or perturbation) procedure: Spin Cycle on the original signals and ASH on the expansion coefficients.

### Example 5.2. Discrimination of geophysical acoustic waveforms.

For the detailed background of this problem, see [18] of this volume. Here, we want to discriminate the waveforms (recorded in a borehole with 256 time samples per waveform) propagated through sandstone layers in the subsurface from the ones through shale layers. We have 201 such “sand waveforms” and 201 “shale waveforms.” We used 10-fold cross validation procedure to compute the misclassification rates. We used the local sine dictionary which is easier to deal with the time information than the wavelet packet dictionaries. We used the relative entropy (6) as a discriminant measure and ASH as the pdf estimator again. In this experiment, we examined the dependence of classification performance to the number of important features  $m$  to compare the results obtained in

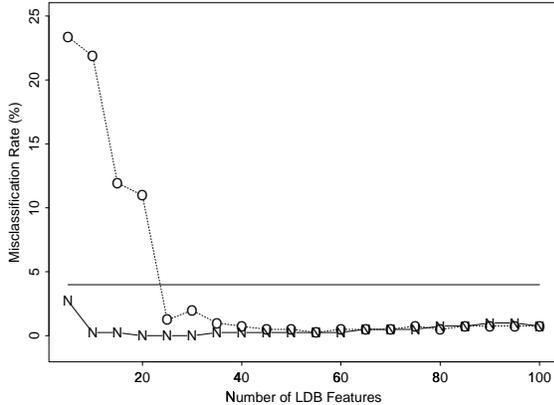


Figure 1: Misclassification rates using LDA as a classifier versus the number of the top LDB features retained. The plots with symbols O and N correspond to the results using the original and the new LDB algorithms, respectively. The constant level line about 4% indicates the performance of the LDA directly applied to the signals represented in the standard coordinate system (of 256 time samples).

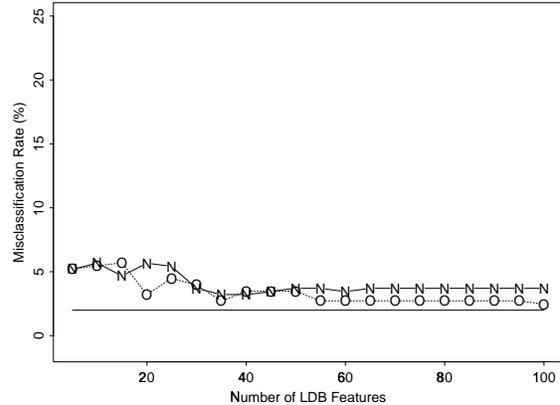


Figure 2: Misclassification rates using CT as a classifier versus the number of the top LDB features retained. The constant level line about 2% indicates the performance of the CT directly applied to the signals represented in the standard coordinate system.

[18]. The results for  $m = k = 5, \dots, 100$  in steps of 5 are summarized in Figures 1 and 2. From these plots, we observe that

- No misclassification occurs with LDA on the top 20, 25, and 30 new LDB vectors.
- These good features are mainly concentrated in P wave components; see also [18].
- Using LDA with less than 40 features, the new LDB outperforms the original LDB. The difference is small for more than 45 features.
- Using CT, the original LDB performs better than the new LDB, but the result on the standard basis is even better.

## 6. CONCLUSION

We described a new LDB algorithm using the “distances” among the estimated pdfs of the projections of input signals onto the basis vectors in the time-frequency dictionaries. Using the probabilistic setting for the new LDB method, the meaning of the original LDB method, which is based on the time-frequency energy distributions of the projections, was clarified. The features derived from the

new LDB vectors can be more sensitive to phase shifts than the original LDB vectors. For the examples we showed, the new LDB method performed better than the original one. We are currently investigating the new LDB method for complex-valued features derived from the local Fourier dictionary, where the new method may have significant advantage over the original one. However, we would like to emphasize that the new algorithm should be considered as an option, not as an absolutely better method than the original one. Depending on the problem, the original LDB method may give sufficient or even better results. In general, one should try both the original and the new LDB methods for her problem at hand.

## 7. REFERENCES

- [1] M. Basseville, *Distance measures for signal processing and pattern recognition*, Signal Processing **18** (1989), no. 4, 349–369.
- [2] L. Breiman, *Bagging predictors*, Tech. Report 421, Dept. of Statistics, Univ. of California, Berkeley, CA 94720, Sep. 1994.
- [3] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and Regression Trees*, Chapman & Hall, Inc., New York, 1993, previously published by Wadsworth & Brooks/Cole in 1984.

- [4] J. B. Buckheit and D. L. Donoho, *Time-frequency tilings which best expose the non-Gaussian behavior of a stochastic process*, Proc. International Symposium on Time-Frequency and Time-Scale Analysis, IEEE, Jun. 18–21, 1996, Paris, France.
- [5] R. R. Coifman and D. Donoho, *Translation-invariant de-noising*, Wavelets and Statistics (A. Antoniadis and G. Oppenheim, eds.), Lecture Notes in Statistics, Springer-Verlag, 1995, pp. 125–150.
- [6] R. R. Coifman and N. Saito, *Constructions of local orthonormal bases for classification and regression*, Comptes Rendus Acad. Sci. Paris, Série I **319** (1994), no. 2, 191–196.
- [7] T. M. Cover, *The best two independent measurements are not the two best*, IEEE Trans. Syst. Man Cybern. **SMC-4** (1974), no. 1, 116–117.
- [8] D. L. Donoho, I. M. Johnstone, G. Kerkycharian, and D. Picard, *Density estimation by wavelet thresholding*, Ann. Statist. **24** (1996), no. 2, 508–539.
- [9] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*, Chapman & Hall, Inc., New York, 1993.
- [10] J. H. Friedman and J. W. Tukey, *A projection pursuit algorithm for exploratory data analysis*, IEEE Trans. Comput. **23** (1974), 881–890.
- [11] R. Guglielmi, *Wavelet Feature Definition and Extraction for Classification and Image Processing*, Ph.D. thesis, Yale University, 1996, In preparation.
- [12] P. J. Huber, *Projection pursuit (with discussion)*, Ann. Statist. **13** (1985), no. 2, 435–525.
- [13] M. C. Jones and R. Sibson, *What is projection pursuit? (with discussion)*, J. R. Statist. Soc. A **150** (1987), no. Part 1, 1–36.
- [14] T. Kohonen, *The self-organizing map*, Proc. IEEE **78** (1990), no. 9, 1464–1480, Invited Paper.
- [15] J. B. Kruskal, *Toward a practical method which helps uncover the structure of a set of multivariate observations by finding the linear transformation which optimizes a new ‘index of condensation’*, Statistical Computation (R. C. Milton and J. A. Nelder, eds.), Academic Press, New York, 1969, pp. 427–440.
- [16] Y. Meyer, *Wavelets: Algorithms and Applications*, SIAM, Philadelphia, PA, 1993, Translated and revised by R. D. Ryan.
- [17] N. Saito, *Local Feature Extraction and Its Applications Using a Library of Bases*, Ph.D. thesis, Dept. of Mathematics, Yale University, New Haven, CT 06520 USA, Dec. 1994, Available via World Wide Web, <http://www.math.yale.edu/pub/wavelets/papers/lfeulb.tar.gz>.
- [18] ———, *Classification of geophysical acoustic waveforms using time-frequency atoms*, ASA Statistical Computing Proceedings, Amer. Statist. Assoc., 1996.
- [19] N. Saito and R. R. Coifman, *Local discriminant bases*, (A. F. Laine and M. A. Unser, eds.), Jul. 1994, Proc. SPIE 2303, pp. 2–14.
- [20] ———, *Local discriminant bases and their applications*, J. Mathematical Imaging and Vision **5** (1995), no. 4, 337–358, Invited paper.
- [21] D. W. Scott, *Multivariate Density Estimation: Theory, Practice, and Visualization*, John Wiley & Sons, New York, 1992.
- [22] F. Warner, *personal communication*, 1996.
- [23] S. Watanabe and T. Kaminuma, *Recent developments of the minimum entropy algorithm*, Proc. Intern. Conf. Pattern Recognition, IEEE, 1988, pp. 536–540.
- [24] M. V. Wickerhauser, *Adapted Wavelet Analysis from Theory to Software*, A K Peters, Ltd., Wellesley, MA, 1994, with diskette.