# QRS MORPHOLOGY REPRESENTATION AND NOISE ESTIMATION USING THE KARHUNEN-LOEVE TRANSFORM

George B. Moody and Roger G. Mark

Massachusetts Institute of Technology, Cambridge, MA, USA

## Summary

We describe robust methods for deriving Karhunen-Loève (KL) basis functions which can be used to represent the QRS complex. Using a five term KL expansion of a 200 millisecond interval which includes the QRS complex and part of the ST segment, we can represent morphology on two simultaneous ECG leads with sufficient fidelity for beat classification. The residual error of the representation is an ideal estimate of the instantaneous noise content of the signal, and permits identification of events for which the morphologic information is unreliable. We have compared the performance of our current KL-based arrhythmia analysis program with its predecessor (which uses a set of time-domain features for morphology representation, but is otherwise identical to the newer program). In evaluations using the MIT-BIH and AHA databases, and a newly developed database containing approximately 2.5 million annotated beats (including over 80,000 PVCs) from 27 long-term ECG recordings, we found that beat classification errors using the KL transform were as little as one-fourth of those for the older program.

## QRS morphology representation

Representation of QRS morphology in terms of a set of numerical parameters is a critical step in automated ECG analysis. If the ECG has been well-sampled (i.e., frequency band-limited and sampled at a rate no less than twice the highest frequency present in the signal), the QRS complex can be reconstructed from samples in the neighborhood of the R-wave peak with sufficient fidelity to permit accurate visual analysis. We define the $n$-dimensional *QRS pattern vector*, $x$, by selecting $n$ baseline-corrected samples surrounding the R-wave peak, each of which is a component of $x$. Such a set of samples thus contains *sufficient* information for analysis; in general, however, its information content far exceeds what is *necessary*. Much of the information encoded into $x$ is redundant, and much of the rest is not relevant to the problem of analyzing the ECG (for example, it may be a faithful representation of noise).

A well-chosen method for data reduction can retain information related to the signal while discarding information related to noise and artifact. For example, bandpass filtering the ECG is equivalent to data reduction in the frequency domain. It is effective because it can remove extraneous components of the signal which are unrelated to cardiac electrical activity. The effectiveness of bandpass filtering for noise rejection is limited, however, since noise is not completely separable from the signal in the frequency domain. The domain in which noise and signal are most separable is that of the Karhunen-Loève transform (KLT). This can be shown to follow directly from the definition of the KLT.

## The Karhunen-Loève transform

The discrete KLT is a rotational transformation of an $n$-dimensional Euclidean pattern space, $E^n$, derived from a model of the probability density function of the pattern vectors, $x$. The probability density function is usually estimated from observations of the distribution of $x$ for a large, representative sample of patterns. If we observe the distribution of $x$ for a large collection of QRS complexes, we find that most of the volume of $E^n$ is devoid of observations. To the extent that the value of a given sample is predicted by those adjacent to it, we find that QRS vectors cluster in small volumes of $E^n$.

Using only second-order statistics, we can model the distribution of $x$ as a hyperellipsoid with principal axes which are defined by the eigenvectors, $e_k$, of the sample covariance matrix

$$C = E\{(x-m)(x-m)'\}$$

where $m$ is the mean pattern vector. The lengths of the axes are given by the corresponding eigenvalues, $\lambda_k$.

Since the $n$ eigenvectors are mutually orthogonal, they span the original pattern space $E^n$. If they are sorted by eigenvalue such that

$$\lambda_k \geq \lambda_{k+1}, \; k=1,2,...,n-1$$

we may expect that, on average, the projection of a pattern vector onto the ordered eigenvectors will have the property that

$$x'e_k \geq x'e_{k+1}$$

We may thus construct $i$-dimensional *feature vectors*, $y_i$, the components of which are given by the first $i$ such projections, $1 \leq i \leq n$; the process of doing so is defined as the discrete KLT. The eigenvectors are the KL basis functions, analogous to the sine functions which are the basis functions of the Fourier transform. The components of $y_i$ are called the KL coefficients, or the principal components, of $x$. For

$i < n$, the KLT has the property that the expected value of the residual error, $r_i = \| \mathbf{y}_i - \mathbf{x} \|$, is minimum among all possible $i$-dimensional linear transformations of the pattern space, for pattern vectors which belong to the distribution modelled by $\mathbf{C}$. Thus, for any desired reduction in dimensionality, the KLT is an optimum linear transformation. Furthermore, for pattern vectors within $E^n$ which do *not* belong to the distribution modelled by $\mathbf{C}$, the expected value of $r_i$ is maximum among all possible linear transformations of the pattern space, for any choice of $i < n$. Provided that the distribution of $\mathbf{x}$ is non-uniform, this last property of the KLT suggests that outliers may be identified by the values of $r_i$ for suitably chosen values of $i$.

### Applying the Karhunen-Loève transform to the ECG

The discussion above has implicitly assumed that it is feasible to model the distribution of all pattern vectors as a single hyperellipsoid. Given the diversity of QRS morphologies observed in the ECG, however, it is not obvious that this assumption is reasonable for QRS pattern vectors. Our preliminary investigations were consistent with the results reported by Nygårds and Sörnmo,[1] who found that the first 11 (of 64) eigenvectors determined from a population of normal QRS complexes were virtually identical to those determined from a population of ectopic QRS complexes. Based on these observations, we adopted the single-hyperellipsoid model for further study.

The prerequisite to applying the KLT to ECG analysis is derivation of the eigenvectors, for which a variety of methods may be employed. It is necessary to obtain a large, representative sample of QRS waveforms (the "training set"). We used the 44 non-paced records of the MIT-BIH Arrhythmia Database[2] for our training set. Approximately 100,000 QRS complexes were identified using an automated QRS detector. Approximately 200 false detections were found by comparison with the reference annotations for the database; these were discarded. (There were also approximately 200 QRS complexes missed by the detector.) The QRS detector, operating on two ECG leads simultaneously, determined a fiducial point for each detected complex by computing the "center of mass" of peaks in the output of the matched filters used for QRS detection. This technique places the fiducial point at or near the major deflection for a monophasic waveform, and midway between the major deflections of a biphasic waveform, thus avoiding discontinuities in fiducial point placement in the context of subtle morphologic variability. For baseline removal, the ECG was digitally high-pass filtered by subtraction of a phase-corrected one-second moving average.

For each QRS complex, we obtained 24 baseline-corrected samples from each of the two ECG leads during a 200-millisecond period which began roughly 60 milliseconds before the fiducial point. (The original signals, digitized at 360 samples per second, were decimated by a factor of three at this time.) These samples defined the components of the QRS pattern vectors which we studied.

Estimation of the eigenvectors is complicated by the presence of noise in the database. Several investigators have

proposed methods for robust estimation of principal components. These include techniques for estimating the covariance matrix using $M$-estimators[3] and projection-pursuit,[4] for estimating the correlation matrix using multivariate trimming,[5] as well as for estimating the eigenvectors directly from the sample distribution.[6] Estola and Jokipii[7] have investigated the robustness of the discrete KLT in this application, and have compared its discriminating power with that of the Fukunaga-Koontz transform, which requires *a priori* knowledge of the cluster means. They emphasize the influence of jitter (suboptimal fiducial point estimation) on the outcome, and suggest that the discriminating power of the KLT may be enhanced if the training set has jitter comparable to that which will be encountered in the ultimate application. In the current study, this condition is satisfied as a consequence of having obtained the pattern vectors in the training set by use of the same QRS detector and fiducial-point location method as is used in the complete KLT-based ECG analysis program.

Our approach makes use of multivariate trimming to obtain a robust sample covariance matrix. For each of the 44 records, we performed an initial clustering and estimated the cluster means. Using 90% of the members of each cluster (those nearest the estimated mean, using a Mahalanobis distance metric), we calculated trimmed cluster means. We repeated this step twice (using the trimmed mean from the previous iteration as the estimate each time); there were essentially no changes in the estimated means between the last two iterations. Using a kernel-approximation method based on approximately 300 estimated cluster means and populations, we obtained a covariance matrix, from which we calculated first-order eigenvectors using standard techniques. We then repeated the entire procedure, reclustering using feature vectors obtained by applying the first-order KLT to the pattern vectors. The second-order eigenvectors obtained at the end of this process differed only slightly from the first-order set. In figure 1, the components of these eigenvectors (the KL basis functions) are plotted as functions of time.

All of the computation described thus far can be performed off-line. In the course of on-line ECG analysis, performing the KLT requires minimal calculation (in our case, 48 integer multiplications and additions per KL coefficient per QRS complex). The motivation for performing the KLT
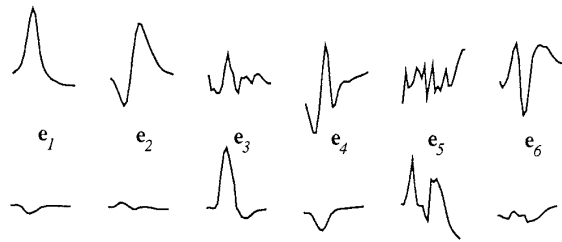


*Figure 1.* The first six QRS eigenvectors (KL basis functions), plotted as functions of time. (Note that each basis function has components related to each of two ECG leads.)

270

in this case is to obtain robust estimates of a few descriptive parameters of the signal (enough to perform reliable clustering) and to differentiate reliable from unreliable morphologic data. For this reason, only a small fraction of the 48 KL coefficients needs to be determined. Figure 2 illustrates a simple case in which two KL coefficients suffice to describe the QRS complexes adequately; as noted below, we find that five KL coefficients are sufficient in general for our ECG analysis program.
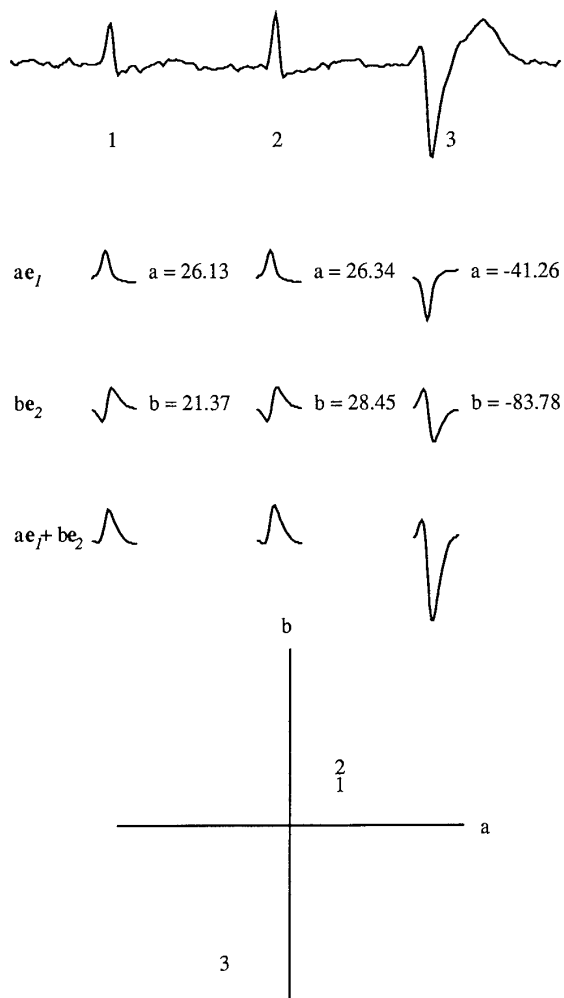


*Figure 2*. Representation of QRS waveforms by summation of KL basis functions. The original ECG is shown at the top of the figure (the second lead is not shown). The first and second KL coefficients (*a* and *b*), and their contributions to the representation of the QRS complexes, are shown in the center of the figure. In the lower part of the figure, the three numbered QRS complexes are mapped onto the *a*–*b* plane, illustrating how the KLT preserves morphologic similarities and differences in preparation for clustering.

## Noise estimation using the KLT

Noise remains the major source of error for state-of-the-art arrhythmia analysis programs, which usually perform well given clean ECGs. We have also applied the KLT to the problem of identifying noise which might lead to errors in analysis. The KLT can represent any input, clean or noisy, given enough coefficients. Using only a few coefficients, however, we can represent most clean inputs much more accurately than most noisy inputs. This observation leads us to use the residual error, $r_i$, of a truncated KLT as an estimate of noise content.

Most other methods for detecting noise in the ECG, including those we have described previously,[8] are based on analysis of the ECG baseline. Such methods cannot succeed in detecting noise bursts which are confined to the QRS complex, which are not uncommon. The present approach has the significant advantage of operating on only that portion of the signal (the QRS complex) in which noise can affect the accuracy of the analysis.

Using a technique we have described previously,[9] we added high-level electrode motion artifact noise to portions of clean ECG records from the MIT-BIH database, choosing the noise level so as to obtain roughly equal numbers of correctly classified beats and errors from our detector. We used ECG records which, without the addition of noise, presented no difficulties to our detector in order to isolate the effects of noise, and we disabled the detector's shut-down logic for this experiment in order to obtain a sufficient number of errors. The correctly-classified beats included essentially all of the beats in the noise-free periods of the input, as well as those beats from the noisy periods which were still recognizable. The errors included both noise-corrupted beats and false QRS detections.

We then measured the residual error, $r_i$, for correctly classified beats and for errors, for various values of $i$. As expected, the residual errors for both groups decrease with increasing $i$, but those for the correctly classified beats (i.e., those from the clean segments of the input signals, and those from the noisy segments which were still recognizable) were lower on average than were the residual errors for the incorrectly classified beats at all values of $i$. The difference between the averages increases with $i$ for small values, as the KL representation of clean QRS complexes improves rapidly. The KL representation of the errors also improves steadily as $i$ increases, however, and the difference between the average errors begins to decrease beyond $i = 5$. Figure 3 illustrates the results of this experiment for a typical record.

## Evaluation

We evaluated our methods using the AHA[10] and MIT-BIH databases, and a newly-developed database containing approximately 2.5 million annotated beats (including over 80,000 PVCs) from 27 long-term ECG recordings. Table 1 summarizes the results. We compared the performance of our current KLT-based analysis program with that of its predecessor. The older program uses time-domain features for QRS morphology representation and noise estimation, but is otherwise identical to the newer one; hence the marked

improvement in PVC detection is attributable entirely to the advantages of the KLT over the time-domain features. (Most of the missed PVCs in the AHA database were in records 7009 and 8007, both of which contain more PVCs than normal beats and present difficulties to the cluster-labelling logic which is common to both programs. Excluding these two records, the PVC sensitivity and positive predictivity are 97.56% and 93.27% respectively using time-domain features, and 97.06% and 95.47% using the KLT.) In tests using the long-term database, we found that the number of missed PVCs was reduced by a factor of four using the KLT as compared to the time-domain features. Less dramatic but still significant improvements were noted for the KLT-based method in analysis of the MIT-BIH and AHA databases.
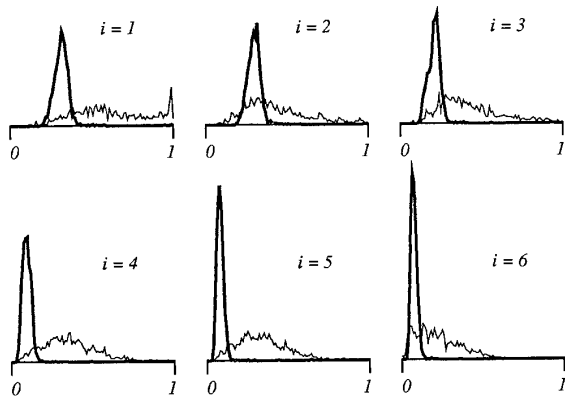


Figure 3. Histograms of normalized residual errors, $r_i/\|x\|$, of the KL morphology representation for correctly classified beats (heavy lines) and incorrectly classified beats (thin lines) of a typical record. The six panels illustrate how the error distributions change with the number of KL coefficients ($i$).

| Database | TDF PVC Se | TDF PVC +P | KLT PVC Se | KLT PVC +P |
|---|---|---|---|---|
| AHA | 91.20% | 93.00% | 91.08% | 94.85% |
| MIT-BIH | 91.90% | 89.46% | 93.12% | 94.76% |
| Long-term | 74.51% | 97.36% | 92.73% | 96.92% |

Table 1. Comparison of the performance of two versions of the same arrhythmia detector. The TDF version uses time-domain features we have described previously[8] and was the subject of optimization efforts between 1981 and 1988. The KLT version is identical except that the components of its feature vectors are KL coefficients, as described in the text. The figures of merit are gross PVC sensitivity (PVC Se) and positive predictivity (PVC +P), as defined by the AAMI.[11]

## Conclusions

We have described a method for deriving a robust QRS morphology description and noise estimation based on the KLT. The QRS morphology description is exceptionally concise (requiring only five integers to describe the QRS complex and a portion of the ST segment on two leads) and is unaffected by moderate amounts of noise. Given noise of sufficiently high amplitude, the signal becomes unrecognizable and the KL coefficients do change; the residual error also increases, however, permitting us to avoid classification error by identifying events in which the morphologic information is unreliable because of the likelihood of noise contamination. In a comparison of the performance of a KLT-based ECG analysis program against a similar program which uses time-domain features, beat classification errors were shown to be significantly lower for the KLT-based method. These results suggest that the KLT is highly effective at retaining the information necessary for an accurate analysis, while rejecting noise and artifact.

References

1. M-E Nygårds and L Sörnmo, "Basis signals representation of QRST waveforms," in Computer-Based Detection of Cardiac Arrhythmias, Linköping University Medical Dissertations No. 151, Linköping, Sweden, 1983. [ISBN 91-7372-658-3]

2. R G Mark and G B Moody, "Evaluation of automated arrhythmia monitors using an annotated ECG database," in Ambulatory Monitoring, ed. C Marchesi, pp. 339-357, Martinus Nijhoff, The Hague, 1984.

3. N A Campbell, "Robust procedures in multivariate analysis I: Robust covariance estimation," Appl Statist, vol. 29, no. 3, pp. 231-237, 1980.

4. G Li and Z Chen, "Projection-pursuit approach to robust dispersion matrices and principal components: primary theory and Monte Carlo," J Am Statist Assoc, vol. 80, no. 391, pp. 759-766, September 1985.

5. S J Devlin, R Gnanadesikan, and J R Kettenring, "Robust estimation of dispersion matrices and principal components," J Am Statist Assoc, vol. 76, no. 374, pp. 354-362, June 1981.

6. K R Gabriel and C L Odoroff, "Resistant lower rank approximation of matrices," in Data Analysis and Informatics, III, ed. E Diday et al., pp. 23-30, Elsevier-North Holland, 1984.

7. K-P Estola and M Jokipii, "Robust features for ECG pattern recognition algorithms," Computers in Cardiology, vol. 14, pp. 253-256, 1987.

8. G B Moody and R G Mark, "Development and evaluation of a 2-lead ECG analysis program," Computers in Cardiology, vol. 9, pp. 39-44, 1982.

9. G B Moody, W K Muldrow, and R G Mark, "A noise stress test for arrhythmia detectors," Computers in Cardiology, vol. 11, pp. 381-384, 1984.

10. R E Hermes and G C Oliver, "Use of the American Heart Association database," in Ambulatory Electrocardiographic Recording, ed. N Wenger, M Mock, I Ringqvist, pp. 165-181, Yearbook Med Pub, Chicago, 1981.

11. Recommended Practice for Testing and Reporting Performance Results of Ventricular Arrhythmia Detection Algorithms, Association for the Advancement of Medical Instrumentation, Arlington, VA, 1986.