

# Journal of Biomedical Optics

[SPIEDigitalLibrary.org/jbo](http://SPIEDigitalLibrary.org/jbo)

## **Photoacoustic spectroscopy of ovarian normal, benign, and malignant tissues: a pilot study**

Sudha D. Kamath  
Satadru Ray  
Krishna K. Mahato

# Photoacoustic spectroscopy of ovarian normal, benign, and malignant tissues: a pilot study

Sudha D. Kamath,<sup>a,b</sup> Satadru Ray,<sup>c</sup> and Krishna K. Mahato<sup>a</sup>

<sup>a</sup>Manipal University, Manipal Life Sciences Centre, Biophysics Unit, Manipal, India

<sup>b</sup>Manipal University, Manipal Institute of Technology, Department of Physics, Manipal, India

<sup>c</sup>Manipal University, Kasturba Medical College, Department of Surgical Oncology, Manipal, India

**Abstract.** Photoacoustic spectra of normal, benign, and malignant ovarian tissues are recorded using 325-nm pulsed laser excitation *in vitro*. A total of 102 (34 normal, 38 benign, and 30 malignant) spectra are obtained from 22 samples belonging to normal, benign, and malignant subjects. Applying multi-algorithm approach, comprised of methods such as, principal component analysis (PCA) based k-nearest neighbor (k-NN) analysis, artificial neural network (ANN) analysis, and support vector machine (SVM) analysis, classification of the data has been carried out. For PCA, first the calibration set is formed by pooling 45 spectra, 15 belonging to each of pathologically certified normal, benign, and malignant samples. PCA is then performed on the data matrix, comprised of the six spectral features extracted from each of 45 calibration samples, and three principal components (PCs) containing maximum diagnostic information are selected. The scores of the selected PCs are used to train the k-NN, ANN, and SVM classifiers. The ANN used is a classical multilayer feed forward network with back propagation algorithm for its training. For k-NN, the Euclidean distance based algorithm is used and for SVM, one-versus-rest multiclass kernel-radial basis function is used. The performance evaluation of the classification results are obtained by calculating statistical parameters like specificity and sensitivity. ANN and k-NN techniques showed identical performance with specificity and sensitivity values of 100 and 86.76%, whereas SVM had these values at 100 and 80.18%, respectively. In order to determine the relative diagnostic performance of the techniques, receiver operating characteristics analysis is also performed. © 2011 Society of Photo-Optical Instrumentation Engineers (SPIE). [DOI: 10.1117/1.3583573]

Keywords: ovarian tissue; photoacoustic spectra; principal component analysis; k-nearest neighbor; artificial neural network; support vector machine; receiver operating characteristics analysis.

Paper 10359RRR received Jun. 25, 2010; revised manuscript received Mar. 24, 2011; accepted for publication Apr. 6, 2011; published online Jun. 1, 2011.

## 1 Introduction

Epithelial ovarian cancer has the highest mortality rate of any of the gynecology cancers and spreads beyond the ovary in 90% of the women diagnosed with ovarian cancer. Ovarian cancer goes undetected in both developed and developing countries because of inadequate technology to detect pre-invasive or early-stage disease. There are approximately 25,600 new cases of ovarian cancer in the United States per year, and there was an estimated 16,000 deaths from ovarian cancer in 2004.<sup>1</sup> High prevalence of late-stage disease and the poor prognosis associated with these later stages are the major factors that give ovarian cancer patients such a dismal prognosis.<sup>1-3</sup>

In recent years, the need for new methods for early detection at a pre-cancer level has stimulated the rapid development of a reliable, objective, and minimally invasive optical tools. The optical (spectroscopic) methods are ideal for this, because in these techniques, spectral changes can be measured with very high sensitivity and they can be used to monitor biochemical alterations, which are the precursors for many diseases including cancer.<sup>4,5</sup> Laser-induced spectroscopic techniques<sup>4</sup> examine different types of light-tissue interactions and noninvasively provide biochemical and morphological information at the molecular, cellular, and tissue levels. They are extremely

sensitive and usually cause no discomfort. Laser-induced spectroscopy techniques,<sup>4-12</sup> such as fluorescence,<sup>4-10</sup> Raman,<sup>11</sup> and photoacoustic<sup>12</sup> (PA), have been widely investigated as nondestructive optical tools for medical diagnostics.<sup>4-13</sup> Pulsed photoacoustic spectroscopy is one of the blossoming fields of science and, in recent years, the technique has become increasingly popular as a tool for biomedical imaging and detection.<sup>13,14</sup> The technique is based on the detection of an acoustic signal induced by photons and is different than the optical techniques detecting optical signals. The propagation of acoustic signals within biological tissues are less susceptible to attenuation and scattering compared to the optical signals and the information contained within them is much less likely to be lost. Because of this advantage, the photoacoustic technique projects itself as a robust and attractive modality for imaging beyond the possible range that exists for all-optical techniques.

The technique of pulsed photoacoustic generation in liquid is well established.<sup>14,15</sup> Scientific applications of this technique have been widely reported in a range of areas including semiconductor research,<sup>16</sup> physical processes in liquids,<sup>14,15</sup> trace gas monitoring,<sup>17</sup> analyte detection in liquids,<sup>15-19</sup> depth-resolved analysis of tissue models,<sup>20</sup> photoacoustic imaging,<sup>21-26</sup> and volumetric analysis of protein.<sup>18,19</sup> A broad range of potential applications have emerged including detection of breast,<sup>22</sup> skin,<sup>20</sup> and oral cancers<sup>27</sup> and vascular applications such as imaging

Address all correspondence to: Krishna K. Mahato, Manipal University, Biophysics Unit, Manipal Life Sciences Centre, Planetarium Complex, Manipal, 576 104, Karnataka, India. Tel: +91-820-2922526; Fax: +91-820-2571919, 2570062; E-mail: kkmahato@gmail.com.

**Table 1** Sample details.

Spectrum number	Sample type	Mean age	Histopathology	Spectroscopy
1 to 15	Normal standard set	50 ± 5.0	Normal	Normal
16 to 30	Malignant standard set	50 ± 20	Ovarian tumor, hemorrhage, cystic, Papillary cystadenoma	Malignant
31 to 45	Benign standard set	40 ± 8.5	Endometriatic cyst	Benign
46 to 64	Normal test set	50 ± 8.5	Normal	Normal
65 to 79	Malignant test set	45 ± 10.5	Papillary cystademinine	Malignant
80 to 102	Benign test set	40 ± 8.5	Benign	Normal, malignant, and benign

superficial blood vessels<sup>28,29</sup> and the characterization of arterial tissues.<sup>29</sup>

In pulsed photoacoustic spectroscopy for biological samples, when a short laser pulse interacts with a tissue, during the laser pulse duration, some of the inherent biomolecules whose concentration depends upon the condition of the tissue having resonance absorption at the excitation absorb the incident photon energy and get excited. The excited biomolecules are then relaxed through nonradiative relaxations releasing the absorbed energy in the form of heat in the absorbing volume of the sample. Subsequently, the thermal expansion of the instantaneously heated sample followed by its contraction due to the periodically applied excitation causes a pressure variation in the irradiated volume. This pressure variation is nothing but the acoustic waves that travel outward through the medium and can be detected at the surface.<sup>30</sup> If the tissue is placed in a cuvette forming the PA cell,<sup>13</sup> the acoustic waves propagating outward from the excitation volume get reflected from the cuvette walls and superpose to form standing waves of different modes depending upon the tissue structure through which they are originated. A few superposed acoustic waves with slight damping will result in complex and long duration time domain photoacoustic transients.<sup>31</sup>

The goal of the present work is to explore whether photoacoustic spectroscopy could eventually be adapted to improve early diagnosis of ovarian neoplasia. With this idea, in our laboratory, we have recorded the photoacoustic spectra of normal as well as different stages of ovarian cancer tissue samples *in vitro* and principal component analysis (PCA) (Refs. 7–10 and 32) based k-nearest neighbor (k-NN),<sup>6,8,9</sup> artificial neural network (ANN),<sup>7,10</sup> and support vector machine (SVM) (Refs. 33 and 34) classification algorithms are developed. This multi-algorithm approach, comprised of computational methods such as ANN, k-NN, and SVM, was used to improve the reliability in classification and hence to determine their relative diagnostic performance.

## 2 Materials and Methods

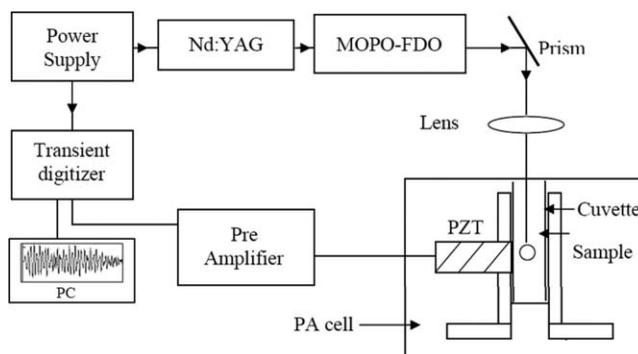
### 2.1 Sample Collection and Handling

Biopsied or surgically resected ovarian normal, benign, and malignant tissue samples of nearly 5×5 to 5×6 mm<sup>2</sup> size were obtained from the Department of Obstetrics and Gynecology, Kasturba Hospital, Manipal University, Manipal, India. Tis-

ues from uninvolved areas from the same subjects are used as healthy controls. A mirror image of each sample was fixed in 10% neutral buffered formalin and was sent for histopathological certification. The samples were kept moist with saline (pH = 7.4), and spectra were recorded within half an hour of tissue removal. The sample details are given in Table 1. According to the histopathology report, samples used in this study were stratified into three main categories: normal ovary (no structures except stroma, epithelium, corpus albicans, and corpus luteum), benign neoplasm (abnormal growth without invasive areas) and endometriosis (growth of both endometrial glands and stroma on the ovary), and malignant (invasion of carcinoma into the ovary, ovarian tumor, hemorrhage, cystic, Papillary cystadenoma).

### 2.2 Experimental Setup

The block diagram of the experimental set up used to record the ovarian tissue photoacoustic spectra is shown in Fig. 1. The samples were excited using 325-nm pulsed laser radiation obtained from an Nd-YAG/MOPO/FDO source (Spectra Physics, Quanta Ray, Model: PRO 230 10, MOPO SL) at 10 Hz repetition of 5 ns pulses with energy 100 to 200 μJ. An indigenously developed PA cell consisting of the “sample holder” and the “PZT mounting” was utilized to generate the PA signal.<sup>12,13</sup> The laser light was focused onto the moist tissue samples kept in a quartz cuvette positioned in the sample holder of the PA cell. The focusing of the laser beam was achieved using a 10 cm focal length convex lens



**Fig. 1** Block diagram of the experimental setup used to record pulsed laser induced photoacoustic spectra.

vertically placed above the cuvette. The acoustic transient signal generated upon laser excitation of the samples was detected with a PZT detector (3 mm diameter and 4 mm thickness, PI-ceramic, Germany) mounted in the photoacoustic cell<sup>12-14</sup> held in contact with the cuvette sidewall. After suitable amplification using a home-made pre-amplifier, the transient signals were recorded in time domain with transient digitizer (EG&G Models 9826 and 9846) personnel computer combination. By changing the incidence position of the laser beam on the sample under study more than two PA transients were recorded from each tissue sample. In the present study, 102 PA transients (34 normal, 38 benign, and 30 malignant) were recorded involving ten normal, 13 benign, and nine malignant tissue samples. All of the PA transients were recorded with maximum care using fixed excitation geometry, as well as with consistent recording procedures.

### 2.3 Fast Fourier Transform of PA spectra

Since the PA cell is essentially a resonant cavity, it is most informative to Fourier transform of the acoustic waves providing the frequency response of the cell and its resonant frequencies for the different types of samples. By examining the signal in the frequency domain one can easily observe the modal frequencies that developed upon excitation of the sample under study. Thus, the acoustic waves (PA signals) originated from different pathological samples, such as normal, benign, and malignant, have slightly different modal frequencies giving rise to slightly different spectral profiles that are directly related to their pathological states. Each pathological sample, thus, produces its own characteristic spectral profile slightly different from one another depending upon the state of the tissue. The spectral features, if extracted from each type of samples, will thus carry significantly different optical properties representing the sample. In the present study, various spectral features extracted from different pathological samples are used for discrimination analysis of the samples. Before extracting spectral features, all of the recorded time domain PA transients were smoothed and Fourier transformed using MATLAB@R6 (Ref. 35) tools, Medfilt1, and FFT (fast Fourier transform), respectively, and the frequency region 0 to  $4 \times 10^5$  Hz of the corresponding FFT pattern was windowed. Subsequently using MATLAB@R6 algorithms, six different features from each of 102 FFT spectra were extracted. These features were mean, median, spectral residual, energy, standard deviation, and maximum intensity, and involving them, a feature space matrix is formed. The process of feature extraction, thus, transformed the frequency domain spectra into a feature space matrix, reducing the number of computations needed for further data analysis. Different spectral features extracted from each of 102 spectra are described below.

- A. Mean: For vectors,  $MEAN(X)$  is the mean value of the elements in  $X$ . Mean of a spectrum is the average intensity over the data points considered in the spectrum.
- B. Median: For vectors,  $MEDIAN(X)$  is the median value of the elements in  $X$ . The median of a spectrum is the intensity at the middle of a group of intensities over the data points considered in the spectrum that have been arranged in order by size.
- C. Standard deviation: For vectors,  $STD(X)$  returns the standard deviation. Standard deviation of a spectrum is the

measure of the dispersion of a set of intensity over the data points considered in the spectrum from their mean. It is the root mean squared deviation.

- D. Energy: Energy of a spectrum is the spectrally integrated intensity over the data points considered in the spectrum.
- E. Spectral residual: A tenth degree polynomial curve is fitted onto normal, benign, and malignant spectra and residual values are noted. This was repeated for all spectra. The norm of residuals is a measure of the goodness of fit, where a smaller value indicates a better fit. The norm is the square root of the sum of the squares of the residuals.
- F. Maximum intensity: It is the maximum intensity of the Fourier transformed frequency domain photoacoustic pattern.

First, the four features, mean, median, standard deviation, and maximum intensity, were obtained using the inbuilt tool, "Data Statistics," the feature Energy using the function 'SUM (X), and the spectral residuals using the Basic Fitting tool of the MATLAB software. Thus, six features extracted from a total of 102 spectra produced an original data matrix of dimension ( $6 \times 102$ ).

### 2.4 Data Analysis

In many of the earlier studies with fluorescence<sup>6,7,10</sup> and Raman spectroscopy,<sup>11</sup> when a number of spectra are recorded from different sites of a tissue sample, the mean of all of the spectra was taken as representative for that sample. In the present study, all of the 102 PA transients recorded are treated as independent data. This, we believe, is a better approach because often a tissue specimen can have normal and neoplastic regions adjacent to each other showing variations in the spectra from site to site.<sup>10,36</sup> The analysis of the spectral data for their classification between different groups has been carried out using PCA based k-NN, ANN, and SVM algorithms.

#### 2.4.1 Principal component analysis

PCA can be used to reduce the dimensionality of a dataset, while still retaining as much of the variability of the dataset as possible.<sup>7,10,32</sup> It is a classical statistical method that transforms attributes of a dataset into a new set of uncorrelated attributes called principal components (PCs). Applying PCA to the feature space matrix, the original data gets transformed into a set of PC scores. The contribution of each PC to the total variance of spectral data is proportional to its eigenvalue. Higher-order PCs often account for less than 1% of the total variation and represent mostly noise.<sup>7,10,32</sup> PCs that have variance more than 1% in the spectral data are considered as informative PCs.

In PCA, first the calibration sets for each class of samples are prepared and an appropriate number of principal components is selected. In the present study, calibration sets are formed using 45 spectra, 15 randomly selected from each group of pathologically certified normal, benign, and malignant samples. PCA is then performed on the data matrix comprised of the six spectral features extracted from each of the 45 calibration samples. The calibration sets are then optimized by cluster analysis removing outliers and an optimum number of factors (PCs) containing maximum diagnostic information

is selected.<sup>7,10,32</sup> Outlier samples usually arise from incorrect measurements, whether it is in the concentration data (i.e., errors in the primary calibration technique) or the spectral data (i.e., sample handling procedures, environmental control such as temperature, humidity, etc). The outlier samples, if present in a calibration set, will tend to “pull” the model in their direction, causing the predicted concentrations of valid samples to be less accurate (or even erroneous) than if the sample was completely eliminated from the training set. The cluster plot for detecting outliers in the calibration model is generally plotted between the PC scores of the first two principal components because the first two principal components represent the plane of best fit through the data. In the present study, when this graph was plotted with arbitrarily chosen spectra of three classes, data points of some spectra were found inconsistent with other samples of the same class in the calibration model and their scores were far from their respective classes. Those spectra with inconsistent scores were then replaced by other spectra of the same class. This process was repeated for a class until spectra with appropriate scores have been identified and found consistent with other scores in the calibration model. Thus, we have optimized the calibration models by including only those spectra that have consistent scores with other samples in the same class and an appropriate number of PCs is identified for each.

After the PCs with diagnostic information are identified, the scores of the selected PCs are utilized and the k-NN, ANN, and SVM algorithms are trained. These trained algorithms are then used to classify unknown spectra. The algorithms for k-NN, ANN, and SVM classification are constructed using only those PCs that have significantly different projection scores for the normal, malignant, and benign spectra.\*\*\*

#### 2.4.2 Classification algorithm: k-nearest neighbor

Nearest neighbor methods provide an important data classification tool for recognizing object classes in the pattern recognition domain.<sup>6,8,9</sup> In this method, an unknown sample is classified to that class having the most “similar” or “nearest” sample point in the training set of data. In the present study, for classification of ovarian data, the single nearest neighbor method has been used. For classification, a prototype sample is computed from the reference (calibration) set and a given test sample is classified as belonging to the class of the closest prototype. That means, when an input pattern is presented to a k-NN classifier, the classifier computes the  $k$  nearest prototypes to it using a distance (Euclidean distances) measure. Then the classifier assigns the class label using a majority vote among the labels of the  $k$  nearest prototypes. This prototype vector is known as the centroid vector. For any test sample, the three spectral distances (i.e., Euclidean distances), corresponding to “normal centroid,” “malignant centroid,” and “benign centroid” are estimated first and are used in the classification. Details of the k-NN algorithm and steps involved in the classification analysis are explained elsewhere.<sup>6,9</sup>

#### 2.4.3 Classification algorithm: artificial neural network

In the present study, ANN with back propagation was used for classification of the ovarian normal, benign, and malignant samples. The back propagation ANN is a potential method for

finding a relationship between different input variables and binomial output variables.<sup>7,10,36,37</sup> The back propagation algorithm consists of fitting the parameters (weights) of the model by a criterion function, usually mean squared error (MSE) or maximum likelihood, using a gradient optimization method. In back propagation artificial neural network, the error (the difference between the predicted outcome and the true outcome) is propagated back from the output to the connection weights in order to adjust the weights in the direction of minimum error. The design of ANN with back propagation algorithm is explained elsewhere.<sup>10</sup> In the classification using ANN, two output neurons are used and the performance goal was fixed at an accuracy of 0.01. The input passes through two processing elements to reach the output. Since the activation function is a sigmoid function, the network (when trained) can represent a nonlinear classifier of any order. The fact that each component of input data reaches each of the output terminals through several parallel paths provides the network considerable flexibility in deciding the nonlinear input-output relationship.<sup>37,38</sup>

#### 2.4.4 Support vector machine analysis

SVM is basically a binary classifier.<sup>33,34</sup> It can be a one-class, two-class, or multiclass SVM.<sup>33,34</sup> In multiclass SVM, it assigns labels to instances by using a support vector machine, where the labels are drawn from a finite set of several elements. This approach reduces the single multiclass problem into multiple binary problems. In the present study, we have used a multiclass SVM one-versus-rest kernel-radial basis function (rbf) based classifier.<sup>33,34</sup> A common approach, in case of SVMs, is the use of recursive feature elimination (RFE) and elimination of the least important features corresponding to the smallest ranking criterion and training the classifier with the remaining features. But, in the present study, we performed MATLAB@R6 based PCA on the training set data matrix comprising six different spectral features extracted from each spectrum as mentioned earlier and an appropriate number of PCs is selected. The scores of the selected PCs containing maximum diagnostic information are then used to train the classifier and the remaining samples are classified using their respective scores with the trained network. Details of SVM design and execution is explained elsewhere.<sup>33,34</sup>

SVM<sup>multiclass</sup> is an implementation of the multiclass support vector machine (SVM).<sup>33,34</sup> Given a training set of instance-label pairs  $(x_i; y_i)$ ,  $i = 1, \dots, l$  where  $x_i \in R^n$  and  $y \in \{1, -1\}^l$ , the SVMs require the solution of the following optimization problem:

$$\min_{w,b,\xi} \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i$$

$$\text{subject to } y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i, \xi_i \geq 0,$$

where  $R$  is the radius of the ball containing the data,  $n$  is the dimension of the input space,  $l$  is the training set size,  $\Phi(x_i)$  is the real valued function before taking threshold,  $\xi$  is the slack variable,  $w$  is the weight vector,  $b$  is the bias, and  $w^T$  is transpose of weight vector. The set of slack variables  $\xi_i$  allow for the class overlap, controlled by the penalty weight  $C > 0$ . This parameter  $C$ , called the regularization parameter, basically controls the trade-off between the largest margin and lowest

number of errors. For,  $C = \infty$ , no class overlap is allowed. In the present study, for the optimal value  $C = 100000$ , the SVM classifier was trained with the training set.

The steps followed in SVM:

- Step 1: Transform data to the format of SVM software
- Step 2: Conduct simple scaling on the data
- Step 3: Consider the rbf kernel:  $K(x_i, x_j) = \exp(-(x_i - x_j)^2/\sigma^2)$
- Step 4: Use cross-validation to find the best parameter  $C$
- Step 5: Use the best parameter  $C$  and to train the whole training set
- Step 6: Test

#### 2.4.5 Receiver operating characteristics analysis

The performances of the classifiers can be visualized through receiver operating characteristic (ROC) graphs and a particular classifier can also be selected with them.<sup>39,40</sup> The ROC graphs are plotted with the true positive rate (TPR) in the y-axis and false positive rate (FPR) in the x-axis.<sup>39,40</sup> Each point on the ROC graph stands for a pair of sensitivity/specificity values corresponding to a particular decision threshold. In case of multiple ROC curves, the area under the curves (AUC) is generally used to evaluate the performance of the curves. For ROC analysis and subsequently for determining the AUC, in the present study, Statistical Package for the Social Sciences (SPSS) 11.0 software<sup>41</sup> was used. The AUC is equivalent to Mann-Whitney U-statistic<sup>39</sup> (nonparametric test of difference between disease/nondisease test results).

### 3 Results and Discussion

Typical smoothed ovarian tissue photoacoustic transients of normal, malignant, and benign classes (left) and the corresponding FFT pattern (right) are shown in Fig. 2. From Fig. 2, noticeable differences among normal, benign, and malignant conditions have been observed. As we know, tissue is basically inhomogeneous in nature containing different chromophores (e.g., water, oxy-hemoglobin, deoxy-hemoglobin, lipid, cytochrome oxidase, and melanin) and fluorophores (tryptophan, collagen, NADH, FAD, etc.) in it at different proportions.<sup>42,43</sup> A small change in its condition alters its biochemical composition leading to changes in its spectral profile. These are essentially the fingerprints that help to determine the condition of the tissue. The PA signal induced by photons depends on the optical energy deposition at the target tissue, which is the product of the tissue optical absorption coefficient and the local light fluence. The optical properties including the absorption coefficient, scattering coefficient, refractive index, and anisotropy factor of the target tissue and the background medium, are all highly dependent on the excitation wavelength.<sup>42</sup> Thus, the complex PA transients originated from normal, benign, and malignant tissues carry different spectral information related to their pathological states.

#### 3.1 Principal Component Analysis

As mentioned in Sec. 2.4, we have six different features extracted from each of 102 photoacoustic spectra (45 calibration

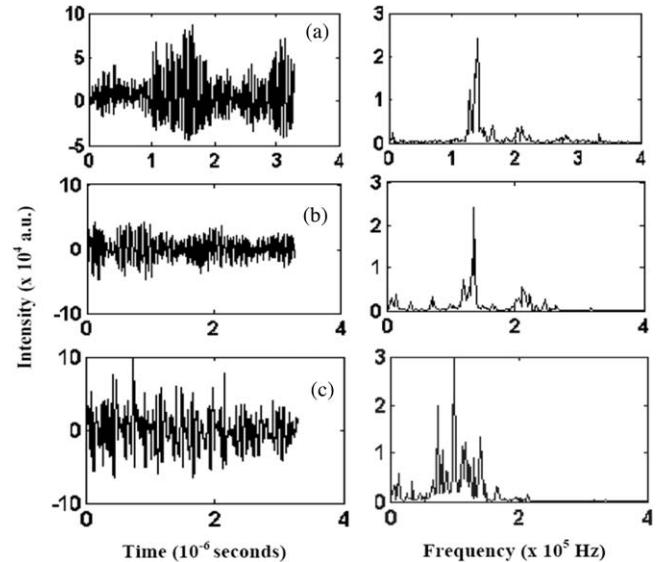
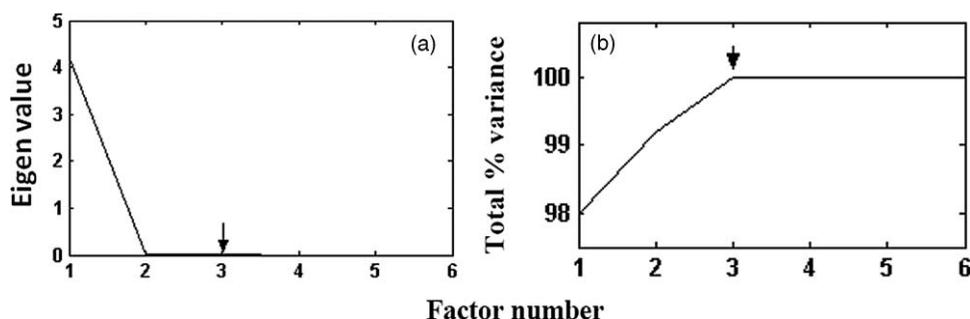


Fig. 2 Typical time domain photoacoustic spectra (left) and the corresponding frequency domain FFT spectra (right) of ovarian tissues. (a) Normal, (b) benign, and (c) malignant.

plus 57 test) used to form a feature space matrix of dimension  $(6 \times 102)$ . On this data matrix, PCA is performed further reducing data dimensionality. In PCA, first, the calibration sets for normal, benign, and malignant classes are formed by randomly selecting 15 pathologically certified spectra from each class. Subsequently, the calibration sets are optimized using cluster analysis removing outliers and an appropriate number of factors (PCs) containing maximum diagnostic information is selected. Figures 3(a) and 3(b) respectively show the eigenvalues for the factors and total percentage variance (i.e., total percentage contribution to the variation spectra with increasing number of factors) of 45 calibration set spectra (15 normal, 15 malignant, and 15 benign). From Figs. 3(a) and 3(b), it is clear that only three factors (PCs) are sufficient to explain the calibration set data. The first PC itself covers over 98.1% of variance and the first three PCs represent 99.92% of the total variance. This shows that the feature vector of length six could be reduced to three components using the PCA technique. As a result, the feature space matrix of dimension  $(6 \times 102)$  was dramatically reduced to dimension  $(3 \times 102)$  making classification computationally more efficient. Thus, we have considered the first three factors for further analysis of the ovarian data.

In PCA, the calibration sets are optimized using cluster analysis removing outliers.<sup>15-17</sup> Figure 4 shows the cluster/scatter plot between the scores of the first two principal components of 45 calibration samples. It is clearly seen from the plot that all of the samples diagnosed as normal, malignant, and benign by pathological examination in the calibration set are clustered in three distinct regions without any overlap, indicating that there are no outliers in these sets. The calibration sets are thus optimized and three informative PCs representing these samples are selected. The scores of these first three factors (PCs) of the 45 (15 normal, 15 malignant, and 15 benign) calibration samples are then used as input feature space for training of k-NN, ANN, and SVM algorithms and the remaining 57 (19 normal,



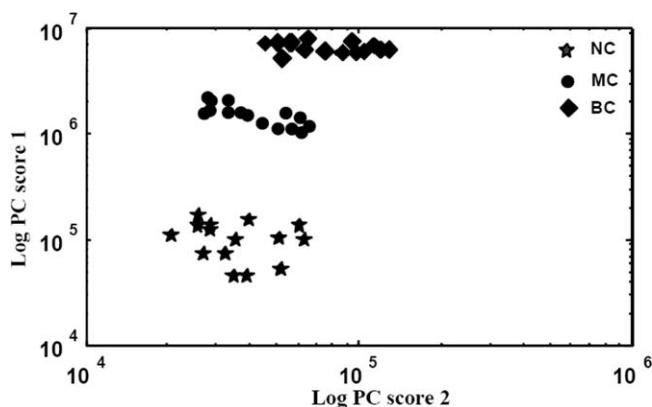
**Fig. 3** Plots of calculated eigenvalues. (a) Total% variance against (b) factors for a PCA decomposition of 45 calibration spectra (15 normal, 15 malignant, and 15 benign).

15 malignant, and 23 benign) samples are tested against them using their respective scores.

### 3.2 *k*-Nearest Neighbor Analysis

For predicting unknown new samples using *k*-NN, first the *k*-NN based programs are executed and then the classification of the unknown samples is carried out as per the algorithm protocol as mentioned in Ref. 6. This technique when used for the 45 calibration samples, all samples are classified to their respective classes. When the 57 test samples (19 normal, 15 malignant, and 23 benign) are tested against the trained algorithm, all normal are classified as normal; two out of 15 malignant are not classified as malignant but as normal; and out of 23 benign, 16 are classified as benign, two are classified as normal, and five are classified as malignant. The two malignant test samples that are classified as normal are the sample numbers 71 and 72 and the two benign test samples which are also classified as normal are the samples 89 and 90. The five benign test samples that are classified as malignant are the samples 95, 96, 97, 98, and 99, respectively. The specificity, sensitivity, and accuracy of this analysis are found to be 100, 86.76, and 91.17%, respectively. The classification results for the calibration and the test samples are shown in Table 2.

Figure 5 shows the plot of Euclidean distance against sample number for the 102 samples (15 + 19 normal, 15 + 15 malignant, and 15 + 23 benign). In Fig. 5, the Euclidean distances for



**Fig. 4** A cluster/scatter plot in log mode between the scores of the first two principal components of the calibration sets spectra. (NC, normal calibration; MC, malignant calibration; BC, benign calibration.)

all of the samples are plotted considering normal centroid as the reference point. The calibration set samples (normal, benign, and malignant) are clustered into three distinct groups without any overlap. The test normal and malignant samples are also clustered along with the corresponding calibration set samples, showing 100% discrimination for these samples. Two of the benign test samples are overlapped with the normal calibration set cluster and five are overlapped with the malignant cluster, showing their tendency toward the normal and malignant. To sum up, only a very small number of malignant/benign samples (two malignant + seven benign) fall outside the general range of malignant/benign species, indicating the probability of malignant/benign samples being in the respective cluster of about 93.33%/81.57% and finding them out of the cluster of about 6.67%/18.43%. The discrimination between different classes of samples is clearly visible in the plot.

### 3.3 Artificial Neural Network Analysis

As mentioned in Sec. 2.4, the ANN algorithm was trained with the input feature vectors formed with the scores of three factors for 15 normal, 15 malignant, and 15 benign samples. Using these feature vectors (15×3 matrix for each case), MATLAB-based ANN programs are executed to train the network and subsequently to predict any new data. The performance goal for the training of ANN was met at 0.001 as shown in Fig. 6 and the convergence was achieved with 11 epochs. The network was given the instruction to show binary digit 1 1 for normal, 1 -1 benign, and -1 -1 for malignant conditions.

When this classification analysis was used on the calibration samples, all 45 samples (15 normal, 15 malignant, and 15 benign) are classified to their respective groups. In this analysis, all of the 19 normal test spectra are classified as normal, two out of 15 malignant test spectra are not classified as malignant but as normal, and 16 out of 23 benign test spectra were classified as benign. Out of the seven remaining benign test spectra, two are classified as normal and five are classified as malignant. The classification results for the calibration and test samples are shown in Table 3. Similar to the *k*-NN analysis, in this case also, the same two malignant test samples 71 and 72 are classified as normal, the same two benign test samples 89 and 90 are classified as normal, and the same five test benign samples 95, 96, 97, 98, and 99 are classified as malignant. This shows that the ANN classification results are exactly the same type as that

**Table 2** Fifteen calibration/test normal/malignant/benign samples tested against normal/malignant/benign calibration/test set.

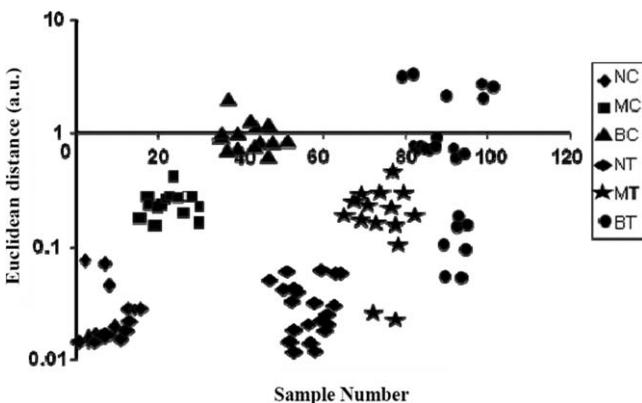
Sample number	Euclidean distances measured from			k-NN classification
	Normal centroid	Malignant centroid	Benign centroid	
1 to 15	0.0139 to 0.0687	0.16 to 0.2566	0.7827 to 0.8753	Normal
16 to 30	0.1475 to 0.2743	0.0011 to 0.0812	0.4427 to 0.7037	Malignant
31 to 45	0.5867 to 2.082	0.358 to 1.8533	0.0291 to 1.2306	Benign
46 to 64	0.0139 to 0.0677	0.161 to 0.2566	0.7853 to 0.8793	Normal
65 to 70	0.1157 to 0.2318	0.0031 to 0.0811	0.6196 to 0.7357	Malignant
71	0.0918	0.1369	0.7596	Normal
72	0.032	0.1967	0.8194	Normal
73 to 79	0.1145 to 0.2264	0.0023 to 0.2293	0.3934 to 0.7369	Malignant
80 to 88	0.334 to 3.1324	0.1311 to 2.9037	0.0758 to 2.281	Benign
89	0.0523	0.1764	0.7991	Normal
90	0.0976	0.1311	0.7531	Normal
91 to 94	0.5897 to 2.082	0.361 to 1.8533	0.1426 to 1.2306	Benign
95 to 99	0.1578 to 0.984	0.0169 to 0.1638	0.6824 to 0.8051	Malignant
100 to 102	1.9223 to 2.652	1.6936 to 2.4233	1.0709 to 1.9006	Benign

of the k-NN results. The specificity, sensitivity, and accuracy of this analysis are also found to be 100, 86.76, and 91.17%, respectively.

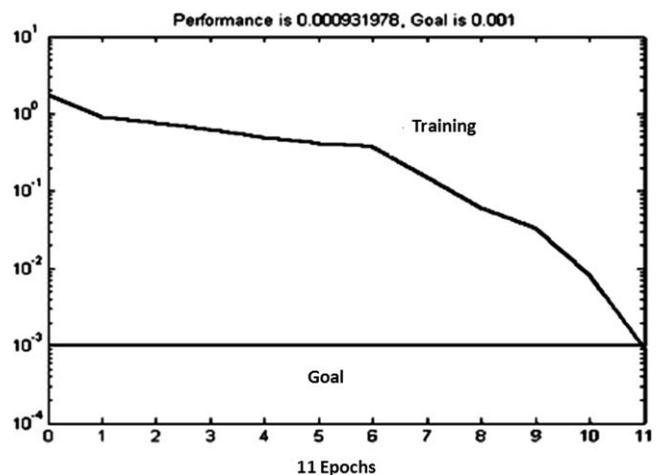
### 3.4 Support Vector Machine Analysis

After constructing the feature space matrix (15×3) for all three calibration sets comprising their PC scores, a

MATLAB-based program was executed and an SVM classifier was trained to predict any new data. The classifier was trained with one-versus-rest (SVM kernel rbf = 10 C = 100,000), where C is the regularization parameter that trades off margin size and training error. The classifier was given the instruction to show binary digit 1 -1 -1 for normal, -1 1 -1 for malignant, and -1 -1 1 for benign conditions.



**Fig. 5** Plot of Euclidean distance versus sample number for the 102 ovarian tissue photoacoustic spectra (15 + 19 normal, 15 + 15 malignant, and 15 + 23 benign). (NC, normal calibration; NT, normal test; MC, malignant calibration; MT, malignant test; BC, benign calibration; BT, benign test.)



**Fig. 6** Training of ANN and its convergence.

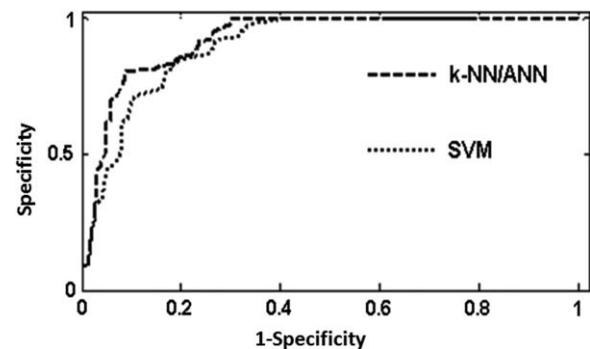
**Table 3** Calibration/test set of normal/malignant/benign (45 + 57 spectra) tested against the trained neural network and results.

Sample number	Desired output		Classifier output		ANN classification
1 to 15	1	1	0.9887 to 1.0000	0.9614 to 1.0000	Normal
16 to 30	-1	-1	-0.9355 to -0.9968	-0.9645 to -0.9999	Malignant
31 to 45	1	-1	0.9584 to 1.0000	-0.9955 to -1.0000	Benign
46 to 64	1	1	0.8575 to 1.0000	0.8859 to 1.0000	Normal
65 to 70	-1	-1	-0.9356 to -1.0000	-0.9583 to -0.9992	Malignant
71	-1	-1	0.0799	0.0928	Normal
72	-1	-1	0.9546	0.9975	Normal
73 to 79	-1	-1	-0.5675 to -0.9959	-0.06315 to -1.0000	Malignant
80 to 88	1	-1	0.9584 to 1.0000	-0.9802 to -1.0000	Benign
89	1	-1	0.9998	0.9975	Normal
90	1	-1	0.9932	0.5466	Normal
91 to 94	1	-1	0.9840 to 0.9998	-0.9198 to -0.9998	Benign
95	1	-1	-0.9860	-0.9872	Malignant
96	1	-1	-0.9736	-0.9891	Malignant
97	1	-1	-0.9195	-0.9943	Malignant
98	1	-1	-1.0000	-0.9522	Malignant
99	1	-1	-0.9736	-0.9802	Malignant
100 to 102	1	-1	0.4632 to 0.9810	-0.9934 to -0.9961	Benign

When this classification technique was used on calibration set samples, all 45 samples (15 normal, 15 malignant, and 15 benign) are classified into their respective groups. In this case, all of the 19 normal test spectra are classified as normal, four out of 15 malignant test spectra are not classified as malignant but as normal, and 14 out of 23 benign test spectra are classified as benign. Out of the nine remaining benign test spectra, six are classified as normal and three are classified as malignant. The calibration and test sample results are given in Table 4. Similar to the k-NN/ANN analysis, in this case also, the same two malignant test spectra 71 and 72 are classified as normal and the same two benign test spectra 89 and 90 are classified as normal. Also, the same three benign test spectra 97, 98, and 99 that are classified as malignant by k-NN/ANN analysis are classified as malignant by SVM. But, the two benign test spectra 95, 96 that are classified as malignant in k-NN/ANN analysis are classified as normal in SVM analysis. In addition to this, two more benign test spectra, 86 and 87 that are classified as benign by the other two classifiers, are classified as normal by SVM. This shows that the SVM classification results are not exactly the same type as that of the k-NN/ANN results. The specificity, sensitivity, and accuracy of the SVM analysis are found to be 100, 80.18, and 90.09%, respectively.

### 3.5 Receiver Operating Characteristic Analysis

After plotting the ROC (Refs. 39 and 40) curves for k-NN/ANN and SVM classifiers, their relative diagnostic performance was determined by measuring the AUC for each curve. The ROC curves plotted with FPR versus TPR for the two classifiers is shown in Fig. 7. The curves are drawn by showing each and every sensitivity/specificity pair resulting from the continuously



**Fig. 7** ROC graph showing relative diagnostic performance of the three discrete classifiers, k-NN, ANN, and SVM analyses.

**Table 4** Calibration/test set of normal/malignant/benign (45 + 57 spectra) tested against the trained SVM classifier and results.

Sample number	Desired output			Classifier output			SVM classification
1 to 15	1	-1	-1	1	-1	-1	Normal
16 to 30	-1	1	-1	-1	1	-1	Malignant
31 to 45	-1	-1	1	-1	-1	1	Benign
46 to 64	1	-1	-1	1	-1	-1	Normal
65 to 70	-1	1	-1	-1	1	-1	Malignant
71	-1	1	-1	1	-1	-1	Normal
72	-1	1	-1	1	-1	-1	Normal
73	-1	1	-1	-1	1	-1	Malignant
74	-1	1	-1	1	-1	-1	Normal
75	-1	1	-1	1	-1	-1	Normal
76 to 79	-1	1	-1	-1	1	-1	Malignant
80 to 85	-1	-1	1	-1	-1	1	Benign
86, 87	-1	-1	1	1	-1	-1	Normal
88	-1	-1	1	-1	-1	1	Benign
89, 90	-1	-1	1	1	-1	-1	Normal
91 to 94	-1	-1	1	-1	-1	1	Benign
95 to 99	-1	-1	1	-1	1	-1	Malignant
100 to 102	-1	-1	1	-1	-1	1	Benign

varying decision threshold over the entire range of results observed in all three analyses. The calculated values of the AUCs for k-NN/ANN and SVM are found to be 0.86 and 0.82, respectively.

#### 4 Conclusion

The pulsed laser induced photoacoustic spectroscopy study conducted on ovarian tissues and the subsequent statistical analysis using PCA-based kNN/ANN/SVM algorithms was motivated by the idea of developing an objective and sensitive photoacoustic-based technique for optical pathology. As a preliminary investigation, the current study reports the discrimination of a limited number of normal, benign, and malignant ovarian tissues *in vitro*. The photoacoustic spectroscopy in combination with PCA-based k-NN/ANN/SVM algorithms has properly and efficiently classified the ovarian tissues suggesting its possible potential application in the field as an alternative or complementary technique to the existing other conventional methods of disease diagnosis. The small time needed to acquire and analyze the photoacoustic spectra together with the high rates of success proves that the technique is very attractive for real time applications. Although, the technique in its current form may not be suitable for such applications, however, with proper instrumentation using a fiber optic probe and a thin film

PZT detector coupled to an endoscope, the technique may be suitable. Of course, for clinical validation of the methodology, further studies with a sufficient number of subjects belonging to normal, benign, and malignant classes as well as on blind samples are extremely essential. Overall, our results indicate that a PCA-based multi-algorithm approach has great promise to classify high dimensional ovarian tissue photoacoustic data and appears to be highly capable to detect ovarian carcinoma, which would be a big improvement in guiding biopsies and diagnosing tissues in different pathological conditions.

#### Acknowledgments

The authors are thankful to the Manipal University, Manipal and to Dr. K. Satyamoothy, Professor and Director, Manipal Life Sciences Centre for providing the necessary facilities to carry out this study. The authors are also thankful to Dr. Rani A. Bhat for providing some of the ovarian tissues used in the study.

#### References

1. M. Brewer, U. Utzinger, E. Silva, D. Gershenson, J. R. C. Bast, M. Follen, and R. R. Kortum, "Fluorescence spectroscopy for *in vivo* characterization of ovarian tissue," *Lasers Surg. Med.* **29**, 128–135 (2001).

2. M. Brewer, U. Utzinger, W. Satterfield, L. Hill, D. Gershenson, R. Bast, J. T. Wharton, R. R. Kortum, and M. Follen, "Biomarker modulation in a nonhuman Rhesus primate model for ovarian cancer chemoprevention," *Cancer Epidemiol. Biomarkers Prev.* **10**, 889–893 (2001).
3. M. A. Brewer, U. Utzinger, J. K. Barton, J. B. Hoying, N. D. Kirkpatrick, W. R. Brands, J. R. Davis, K. Hunt, S. J. Stevens, and A. F. Gmitro, "Imaging of the ovary," *Technol. Cancer Res. Treat.* **3**(6), 617–627 (2004).
4. N. Ramanujam, "Fluorescence spectroscopy of neoplastic and non-neoplastic tissues," *Neoplasia* **2**, 89–117 (2000).
5. A. Pradhan, P. Pal, G. Durocher, L. Villeneuve, A. Balassy, F. Babai, L. Gaboury, and L. Blanchard, "Steady state and time resolved fluorescence properties of metastatic and non-metastatic malignant cells from different species," *J. Photochem. Photobiol., B* **31**, 101–112 (1995).
6. S. D. Kamath and K. K. Mahato, "Optical pathology using oral tissue fluorescence spectra: classification by principal component analysis (PCA) and k-means nearest neighbour (k-NN) analysis," *J. Biomed. Opt.* **12**(1), 014028 (2007).
7. S. D. Kamath, C. D'souza, S. Mathew, S. George, S. Chadangil, and K. K. Mahato, "A pilot study on colonic mucosal tissues by fluorescence spectroscopy technique: discrimination by principal component analysis (PCA) and artificial neural network (ANN) analysis," *J. Chemom.* **22**, 408–418 (2008).
8. S. D. Kamath, R. A. Bhat, S. Ray, and K. K. Mahato, "Autofluorescence of ovarian normal, benign, and malignant tissues: a pilot study," *Photomed. Laser Surg.* **27**(2), 325–335 (2009).
9. S. D. Kamath and K. K. Mahato, "Principal component analysis (PCA) based k nearest neighbour (k-NN) analysis of colonic mucosal tissue fluorescence spectra," *Photomed. Laser Surg.* **27**(4), 659–668 (2009).
10. G. S. Nayak, S. D. Kamath, K. M. Pai, A. Sarkar, S. Ray, J. Kurien, L. D'Almeida, B. R. Krishnanand, S. Chadangil, V. B. Kartha, and K. K. Mahato, "Principal component analysis (PCA) and artificial neural network (ANN) analysis of oral tissue fluorescence spectra: classification of normal, pre-malignant, and malignant pathological conditions," *Biopolymers* **82**, 152–166 (2006).
11. C. M. Krishna, G. D. Sockalingum, J. Kurien, L. Rao, L. L. Venteo, M. Pluot, M. Manfait, and V. B. Kartha, "Micro-Raman spectroscopy for optical pathology of oral squamous cell carcinoma," *Appl. Spectrosc.* **58**, 107–114 (2004).
12. R. A. Bhat, S. D. Kamath, K. K. Mahato, S. Ray, and V. B. Kartha, "Photoacoustic spectroscopic studies of ovarian tissue in different pathological conditions: classification using cluster analysis," *Proc. Amer. Assoc. Cancer Res.*, **47**, No. 3584 (2006).
13. S. D. Kamath, V. B. Kartha, and K. K. Mahato, "Dynamics of L-tryptophan in aqueous solution by simultaneous laser induced fluorescence (LIF) and photoacoustic spectroscopy (PAS)," *Spectrochim. Acta Part A* **70**(1), 187–194 (2007).
14. C. K. N. Patel and A. C. Tam, "Pulsed optoacoustic spectroscopy condensed matter," *Rev. Mod. Phys.* **53**, 517–550 (1981).
15. M. W. Sigrist, "Laser generation of acoustic waves in liquids and gases," *J. Appl. Phys.* **60**, R83–R121 (1986).
16. A. Rosencwaig and A. Gersho, "Theory of the photoacoustic effect with solids," *J. Appl. Phys.* **47**, 64–69 (1976).
17. P. Hess, "Photoacoustic, photothermal and photochemical processes in gases," in *Topics in Current Physics*, Vol. 46, Springer-Verlag, Berlin (1989).
18. S. E. Braslavsky, "Photoacoustic and photothermal methods applied to the study of radiationless deactivation processes in biological system and in substances of biological interest," *Photochem. Photobiol.* **43**, 667–679 (1986).
19. T. Schmid, "Photoacoustic spectroscopy for process analysis," *Anal. Bioanal. Chem.* **384**, 1071–1086 (2006).
20. G. Puccetti, F. Lahjomri, and R. M. Leblanc, "Pulsed photoacoustic spectroscopy applied to the diffusion of sunscreen chromophores in human skin: the weakly absorbent regime," *J. Photochem. Photobiol. B* **39**, 110–120 (1997).
21. E. Zhang, J. Laufer, and P. Beard, "Backward-mode multiwavelength photoacoustic scanner using a planar Fabry–Perot polymer film ultrasound sensor for high-resolution three-dimensional imaging of biological tissues," *Appl. Opt.* **47**(4), 561–577 (2008).
22. S. Manohar, S. E. Vaartjes, C. G. Johan, V. Hespden, J. M. Klaase, M. E. Frank, W. Steenbergen, and T. G. van Leeuwen, "Initial results of in vivo non-invasive cancer imaging in the human breast using near-infrared photoacoustics," *Opt. Express* **15**(19), 12277–12285 (2007).
23. L. Xiang, D. Xing, H. Gu, D. Yang, S. Yang, L. Zeng, and W. R. Chen, "Real-time optoacoustic monitoring of vascular damage during photodynamic therapy treatment of tumor," *J. Biomed. Opt.* **12**, 014001 (2007).
24. D. W. Yang, D. Xing, Y. Tan, H. Gu, and S. Yang, "Integrative prototype B-scan photoacoustic tomography system based on a novel hybridized scanning head," *Appl. Phys. Lett.* **88**, 174101 (2006).
25. B. T. Cox, S. R. Arridge, and P. C. Beard, "Quantitative photoacoustic image reconstruction for molecular imaging," *Proc. SPIE* **6086**, 60861M (2006).
26. L. Xiang and F. Zhou, "Photoacoustic imaging application in tumor diagnosis and monitoring," *Key Eng. Mater.* **364–366**, 1100–1104 (2008).
27. A. A. Oraevsky, A. A. Karabutov, E. B. Savateeva, B. Bell, M. Motamedi, S. L. Thomsen, and P. Pasricha, "Opto-acoustic imaging of oral cancer: feasibility studies in hamster model of squamous cell carcinoma," *Proc. SPIE* **3597**, 385–396 (1999).
28. J. Laufer, D. Delpy, C. Elwell, and P. Beard, "Quantitative spatially resolved measurement of tissue chromophore concentrations using photoacoustic spectroscopy: application to the measurement of blood oxygenation and hemoglobin concentration," *Phys. Med. Biol.* **52**, 141–168 (2007).
29. P. C. Beard and T. N. Mills, "Characterization of post mortem arterial tissue using time-resolved photoacoustic spectroscopy at 436nm, 461nm and 532nm," *Phys. Med. Biol.* **42**(1), 177–198 (1997).
30. Y. Su, F. Zhang, K. Xu, J. Yao, and R. K. Wang, "A photoacoustic tomography system for imaging of biological tissues," *J. Phys. D: Appl. Phys.* **38**, 2640–2644 (2005).
31. D. A. Schurig, G. L. Klunder, M. A. Shannon, and R. E. Russo, "Signal analysis of transients in pulsed photoacoustic spectroscopy," *Rev. Sci. Instrum.* **64**(2), 363–373 (1993).
32. I. T. Jolliffe, "Principal component analysis," Springer-Verlag, New York (1986).
33. S. S. Keerthi and C. J. Lin, "Asymptotic behaviors of support vector machines with Gaussian kernel," *Neural Comput.* **15**(7), 1667–1689 (2003).
34. C. C. Chang and C. J. Lin, "Training  $\nu$ -support vector classifiers: theory and algorithms," *Neural Computation* **13**(9), 2119–2147 (2001).
35. J. W. Palm III, "Introduction to MATLAB 6 for engineers," McGraw-Hill, New York (1994).
36. M. F. Mitchell, "Accuracy of colposcopy," *Consult. Obstetrics Gynecol.* **6**, 70–73 (1994).
37. J. M. Zurada, *Introduction to Artificial Neural Systems*, 3rd ed., Jaico Publishing House, Mumbai, India (2002).
38. S. Hykin, *Neural Networks. A Comprehensive Foundation*, 2nd ed., Pearson Education, Singapore (2001).
39. T. Fawcett, "An introduction to ROC analysis," *Pattern Recogn. Lett.* **27**, 861–874 (2006).
40. J. A. Hanley and B. J. McNeil, "The meaning and use of the area under a receiver operating characteristic (ROC) curve," *Radiology* **143**, 29–36 (1982).
41. V. Gupta, "SPSS for-beginners" 2nd edition, VJ Books Inc. (2000); <http://www.brothersoft.com/spss-11-for-beginners-2nd-edition-34060.html>.
42. J. R. Rajian, P. L. Carson, and X. Wang, "Quantitative photoacoustic measurement of tissue optical absorption spectrum aided by an optical contrast agent," *Opt. Express* **17**(6), 4879–4889 (2009).
43. G. Wagnieres, W. Star, and B. C. Wilson, "In vivo fluorescence spectroscopy and imaging for oncological applications," *Photochem. Photobiol.* **68**(5), 603–632 (1998).