

Representation and Classification of Microstructures using Statistical Learning Techniques

Veeraraghavan Sundararaghavan* and Nicholas Zabaras*

*Materials Process Design and Control Laboratory, Sibley School of Mechanical and Aerospace Engineering, 188 Frank H. T. Rhodes Hall, Cornell University Ithaca, NY 14853-3801, USA

Abstract. Microstructural information is fundamental to determining the critical properties of today's high performance materials. Hence, there is a need for a representation that can quantify all the microstructural elements through the analysis of digitized images. This paper addresses representation through the creation of a dynamic microstructure library. The paper focuses on the application of machine learning theory for the creation of a library that is trained by experimentally or computationally obtained microstructure snapshots. Support vector machines (SVM) are used to classify microstructure snapshots based on its features into various classes. An incremental-principal component analysis (PCA) method is employed within image classes to constantly update the microstructure basis and numerically quantify the microstructural features.

MOTIVATION

Models that attempt to predict properties of materials rely on statistical descriptions of material microstructure. Descriptors like grain sizes, elongations and orientations that describe the microstructure morphology belong to the class of lower order geometrical descriptors. Attempts to reconstruct microstructures based on lower order descriptors [1] result in non-uniqueness which can be attributed to the absence of complete morphological information in such lower order measures. This paper focuses on the completeness aspect of the description of planar images of single phase polyhedral microstructures. We employ a classification technique in conjunction with principal component analysis for the quantification of microstructure images.

Jenkins *et al.* [2] and Tojima *et al.* [3] used classification methods for describing material microstructure. However, the descriptors used for the description are of lower order and hence incomplete. The automated classification structure proposed in this paper employs lower order descriptors like grain sizes and shape features for classification but complete description is obtained by describing the microstructure through the use of an image basis within a class library. In this approach, feature component vectors representing independent patterns extracted from the various classes of single phase polyhedral microstructures are used to train a system. Initial training classes consisting of ensembles of single phase polyhedral microstructure were constructed using a Monte Carlo algorithm for grain growth. The library does not store any microstructure images. Every class

consists of a reduced basis which can effectively describe new images. This basis evolves when new images are added to the classifier and the information content of the class improves. New microstructures can be completely reconstructed through a linear combinations of the basis through a set of coefficients, which are used for quantitatively representing the microstructure. A modified form of principal component analysis (PCA) [4] extensively applied in face recognition and vision applications has been employed for the creation of the basis. This method called incremental PCA [5] technique dynamically updates the basis within each class whenever new images are added to the library.

FEATURE EXTRACTION

An automated image analysis scheme is adopted for feature extraction. Raw images are initially modified so that all images in the microstructure library have the same orientation and magnification. Fig. 1 shows the important preprocessing steps such as image alignment, scaling and subsequent steps that involve sharpening the image through illumination equalization, edge enhancement, and finally grain boundary detection. The boundary image is then used in the size and shape parameter identification steps.

The following three feature vectors were extracted from the input microstructure image (Fig. 2(a)):

1. Histogram of the intercept length distribution using

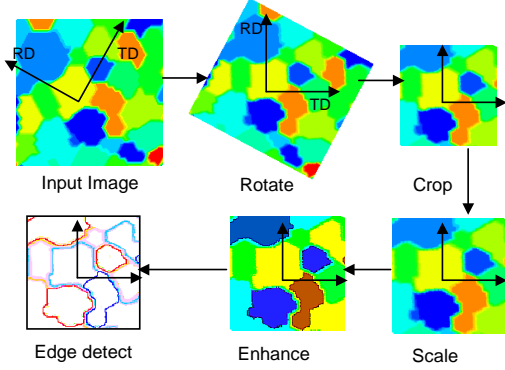


FIGURE 1. Microstructure image preprocessing operations.

Heyn's intercept technique (mean intercept length versus number of test lines possessing the mean intercept length, Fig. 2(b)) [6].

2. The rose of intersections (Fig. 2(c)) [7].
3. Color histogram.

PRINCIPAL COMPONENT ANALYSIS

Principal component analysis is a technique to obtain low-dimensional representation of a large set of data. Using a set of large dimensional data called the 'snapshots', the method decomposes the data into an optimal orthonormal basis. Few basis vectors selected in the order of importance can be used for the representation of the high dimensional data sets. This method is well suited to the representation of images. Fig. 3 shows sample microstructure images (I_i) used as an example to demonstrate PCA.

Let N different microstructure planar images (I_i), each of size n pixels by n pixels are to be represented. The images are converted into N vectors (X_i) and the average is computed as $\mu = \frac{1}{N} \sum_{i=1}^N X_i$. The average image (μ) is then subtracted from all the image vectors as $X_i \leftarrow X_i - \mu$, for $i = 1, \dots, N$. The eigenvectors \vec{U}_k of the $n \times n$ covariance matrix $C = \frac{1}{N} \sum_{i=1}^N X_i X_i^T$ satisfying

$$CU_k = \lambda_k U_k, \quad k = 1, \dots, N \quad (1)$$

with the eigenvalues λ form the best basis for the images. Even though the above method calculates the best uncorrelated basis, it is computationally intensive since it involves the calculation of a correlation matrix of very large dimensionality. The work around is the so called 'method of snapshots'. Here, the property that the eigenvectors U_k are the unique linear combinations of the microstructure images (X_i) is used and the eigenvectors can

TABLE 1. Coefficients of the input images in the eigen basis ($\times 0.001$)

0.0125	1.3142	-4.23	4.5429	-1.6396
-0.8406	0.8463	-3.0232	0.3424	2.6752
3.943	-4.2162	-0.6817	-0.9718	1.9268
1.1796	-1.3354	-2.8401	6.2064	-3.2106
5.8294	5.2287	-3.7972	-3.6095	-3.6515

thus be written as

$$U_k = \sum_{j=1}^N \alpha_{jk} X_j \quad k = 1, \dots, N \quad (2)$$

Let us define C^* as $X_i^T X_j$, $i, j = 1, \dots, N$, and let the vector $E_k = \alpha_{ik}$, $i = 1, \dots, N$, denote the coefficients of the eigenvector U_k in the basis of the snapshots. Then the original eigenvalue problem Equation (1) is equivalent to the eigenvalue problem,

$$C^* E_k = \lambda_k^* E_k \quad (3)$$

A $N \times N$ matrix, $X_i^T X_j$ is constructed and the vectors E_k , $k = 1, \dots, N$ are found from the solution of the above eigenvalue problem. The N eigenvectors U_k are subsequently found using Equation (2). These vectors form the so called 'eigenfaces', U . The eigenfaces for the microstructures in Fig. 3 are shown in Fig. 4. The eigenface vectors are normalized and stored in the material library.

Once the eigen-basis for the set of microstructure in the class is identified, any new image corresponding to that class can be represented by transforming the image into the eigenface components by a projection operation. The coefficients (ω_k) of the new image (Γ) in the normalized eigen-basis is given by

$$\omega_k = U_k^T (\Gamma - \mu) \quad (4)$$

The coefficients (ω_k) form a vector $\Omega = [\omega_1, \dots, \omega_N]^T$ that is used as a reduced representation for the microstructure. The matrix of coefficients of the input images $[\Omega_1, \dots, \Omega_N]$ is denoted by A , the representation matrix. The representation matrix for the input images in the eigen-basis is listed in Table 1.

SUPPORT VECTOR MACHINES

The aim of classification is to group similar microstructures within a class where PCA analysis may be carried out. The image classification problem has been solved using an statistical learning algorithm called Support vector machines [8]. The classification involves prior training with features from known microstructure classes. The training involves finding the optimal hyper-plane such that the error for unseen test microstructure

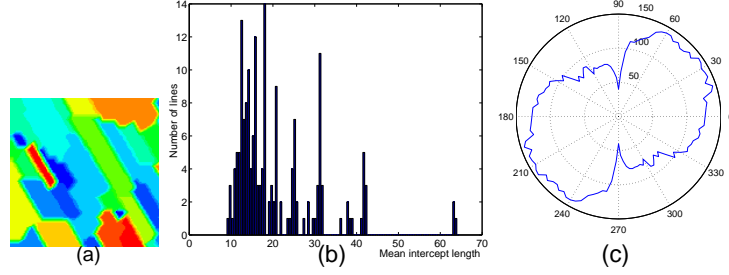


FIGURE 2. (a) The input microstructure image (b) Histogram of mean intercept length versus number of lines obtained using Heyn intercept technique (c) Rose of intersections.



FIGURE 3. Sample microstructure images to be represented using PCA.

images is minimized. Each instance in the training set consists of class labels and several attributes extracted from the microstructure image. The goal of SVM is to produce a model which predicts the class of the data set given in the form of features.

Given a training set of instance-label pairs (x_i, y_i) , $i = 1, \dots, n$, where $x_i \in \mathcal{R}^n$ and $y \in \{1, -1\}$, the support vector machines non-linearly maps the data \vec{x} to a higher dimensional feature space F as $\vec{z} = \phi(\vec{x})$. The function $\phi(\vec{x})$ is defined by a positive definite kernel, $K(\vec{x}, \vec{x}')$, specifying an inner product in the feature space, $\phi(\vec{x}) \cdot \phi(\vec{x}') = K(\vec{x}, \vec{x}')$. The kernel employed for microstructure classification is the linear kernel, $K(\vec{x}, \vec{x}') = \vec{x} \cdot \vec{x}'$. If the data is linearly separable in F , a decision function $D(\vec{x}) = \vec{w} \cdot \phi(\vec{x}) + b$ where \vec{w} is an n -dimensional vector and b is a scalar can be determined such that,

$$y_i D(x_i) \geq 1, \quad i = 1, \dots, n \quad (5)$$

The distance from the separating hyperplane $D(\vec{x}) = 0$ and the training datum nearest to the hyperplane is called the margin. The hyperplane with the highest margin is called the optimal hyperplane. Fig. 5 shows a realization of a binary classifier through the creation of a hyperplane that maximizes the margin between the two examples. Vector \vec{w} for the optimal hyperplane is obtained by maximizing the margin which becomes equivalent to minimizing $\|\vec{w}\|$. The solution to this is given by $\vec{w} = \sum_{i=1}^n \alpha_i y_i \phi(\vec{x}_i)$ for $\alpha_i \geq 0$. The problem of determining the α_i 's is posed as a quadratic programming problem of maximizing,

$$W(\vec{\alpha}) = \sum_{i=1}^n \alpha_i$$

$$- \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j K(\vec{x}_i, \vec{x}_j), \quad (6)$$

in the positive quadrant $\alpha_i \geq 0$, $i = 1, \dots, n$, subject to the constraint, $\sum_{i=1}^n \alpha_i y_i = 0$. The support vectors are the points for which $\alpha_i > 0$ satisfying Equation (5) with equality. Given a new set of data x , the decision function can be written as,

$$f(x) = \text{sgn} \left(\sum_{i=1}^n \alpha_i y_i K(\vec{x}_i, \vec{x}) + b \right) \quad (7)$$

When the data is non-separable in the higher dimensional feature space, slack variables $\xi_i \geq 0$ are introduced such that $y_i D(x_i) \geq 1 - \xi_i$, $i = 1, \dots, n$ to allow the possibility of samples to violate Equation (5). The idea is to maximize the margin and minimize the training error (represented by the slack variables) simultaneously. The generalized optimal separating hyperplane is then the minimization of $\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i$ where the purpose of C in the second term is to control the number of misclassified points.

Microstructure classification involves operating on more than one feature parameter over several classes of images. Such problems are called Multi-class problems. The one-against-one method has been employed for this classification problem [9]. This method constructs $k(k-1)/2$ classifiers where each classifier is trained on data from two classes. If between two classes i and j , a given data set is classified to class i , then the vote for class i is incremented by one. Given a data point, $k(k-1)/2$ classifications are performed and the data is classified to the class which gets the maximum votes. In case that two classes have identical votes (the data set falls in the indecision region) the class with a smaller index is selected.



FIGURE 4. Images of the eigen-basis: the eigenfaces.

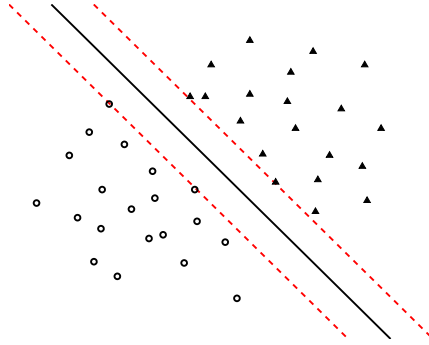


FIGURE 5. Support vectors (shown as dotted lines) used for binary classification.

IMPLEMENTATION

Using the features of representative classes of images as training data sets, the support vector machines are used to group new unknown test examples to these classes. Fig. 6 describes the class structure of the library. The first tier of the library is based on the shape features of microstructures. Further classification of grain shapes on the basis of grain sizes is indicated in the figure. Using a hierarchical classification structure, the user will have the added flexibility to branch off new classes based on additional problem specific microstructural features.

To test the classification scheme employed, 375 raw images of microstructures were created using the Monte Carlo grain growth model. Each image was sized to 128×128 pixels with 256 gray levels and subdivided into 11 classes of microstructures based on grain shapes. The rose of intersections of microstructure images were normalized in the range $[-1, 1]$ and were used as feature vectors. About ten percent of randomly selected images were used for training and the rest were used to test the accuracy of classification. It is ensured that the training and the test sets do not overlap. The classifier was repeatedly trained using random sets of 40 images and the classifier accuracy was checked using the rest 335 images as test examples. On the average, the multi-class support vector machine classification scheme gave an accuracy of 92.53 percent, with the lowest accuracy of 87.76 percent and highest accuracy of 95.82 percent for 30 such experiments. The number of training examples were increased to 100 randomly picked images and the rest 275

images were provided as training examples. An average accuracy of 95.80 percent is realized.

Principal component analysis works dynamically with the classification and the basis is updated with inclusion of new images. Fig. 7 shows images within the class reconstructed using different fractions of eigenvectors in the basis. It can be seen that good quality reconstructions are possible even when 40 percent of the eigen-basis is discarded.

Whenever a new image (Γ_k) of dimension ($L_1 \times L_2$) is added to a class of images, the existing basis (U_k) is used to reconstruct the microstructure. The efficiency of reconstruction is tested using a measure (d) based on the test image (Γ_k) and the reconstructed microstructure (Γ_k^R) given by,

$$d = \frac{\|\Gamma_k - \Gamma_k^R\|_F}{\sqrt{L_1 L_2}} \quad (8)$$

The test and the reconstructed microstructures were normalized within limits of $[0, 1]$ and a distance threshold of $d > 0.1$ was employed as the criteria for updating the basis.

The PCA technique explained earlier is performed in the batch mode and requires the calculation of the correlation matrix and hence, it requires redoing the entire analysis if a new image is included in the basis. In addition, since all images are processed simultaneously at each step, the method requires the storage of all previously coded images. This is inadmissible in a dynamic library where images are processed sequentially. We employ an incremental PCA technique proposed by Skočaj and Leonardis [5] for a dynamic update of the eigenfaces. In this method, the microstructure images are discarded after the PCA update and only the coefficients of the images are stored along with the associated eigen-basis. This provides a framework for online representation of microstructure images. The Euclidean measure Equation (8) is used as the criterion for the update of the basis. The initial eigen basis for this algorithm is obtained by applying the batch PCA on a small set of images in each class of the microstructure library. The representation format of an image in the dynamic material library is shown in Table 2. The representation coefficients are generated using the Eigen basis of images in the combined shape and size class.

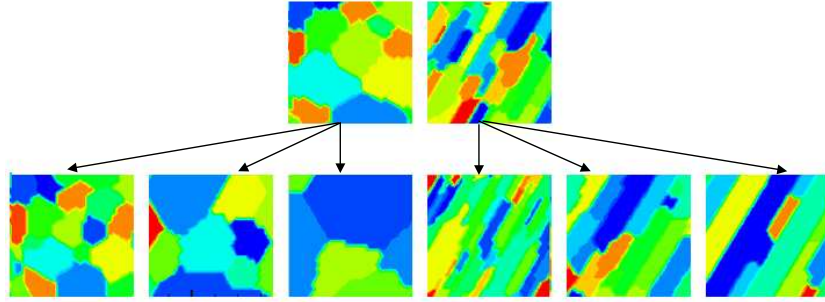


FIGURE 6. Classification scheme for microstructures.

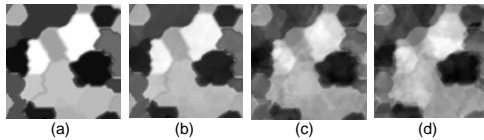


FIGURE 7. Microstructure reconstructed using PCA with: (a) 100 percent of the basis; (b) 80 percent of the basis; (c) 60 percent of the basis; and (d) 40 percent of the basis.

TABLE 2. Representation format for a microstructure in the library

Date: 1/12 02:23PM, Basis updated
Shape Class: 3, (Oriented 40 degrees, elongated)
Size Class : 1, (Small grains)
Coefficients in the basis: [2.42, 12.35, -4.14, 1.95]

CONCLUSIONS

A classification based representation scheme for single phase polyhedral microstructures based on a dynamic microstructure library has been demonstrated. Support vector machines, a statistical learning algorithm, has been used to classify microstructure images on the basis of extracted grain shape and size features. Each class of microstructures generated using support vector machines holds a basis for representing microstructure images. The basis dynamically updates with the arrival of new images within a class. Complete microstructure representation is achieved through a set of coefficients in the generated basis.

ACKNOWLEDGMENTS

The work presented here was funded by the Computational Mathematics program of the Air Force Office of Scientific Research (grants F49620-00-1-0373 and FA9550-04-1-0070) and by the National Science Foundation (grant DMI-0113295). This research was con-

ducted using the resources of the Cornell Theory Center, which receives funding from Cornell University, New York State, federal agencies, and corporate partners.

REFERENCES

1. C.L.Y. Yeong and S. Torquato, *Physical Review E* **57(1)**,495-506 (1998).
2. B.M. Jenkins, R.R. Lovel and J.A. Thurlby, "Quantification of iron ore sinter structure by optical image analysis" in *Metallography: Past, Present and Future (75th Anniversary Volume)*, edited by G.F. Vander Voort et al., ASTM STP 1165, Philadelphia, 1993, pp. 292-299.
3. M. Tojima, T. Suzuki and F. Kobayashi, "Classification of graphite shapes in cast irons using neural networks" in *COMP '93 International Conference on Computer-Assisted Materials Design and Process Simulation Proceedings*, Iron and Steel Inst., Tokyo, Japan, 1993, pp. 463-470.
4. M. Turk and A. Pentland, *J Cognitive Neurosci* **3(1)**,71-86 (1991).
5. D. Skočaj and A. Leonardis, "Incremental approach to robust learning of eigenspaces" in *26th Workshop of the Austrian Association for Pattern Recognition (ÖAGM/AAPR)*, edited by F. Leberl and F. Fraundorfer, Graz (Austria), 2002, pp. 71-78.
6. G.F. Vander Voort, "Examination of some grain size measurement problems" in *Metallography: Past, Present and Future (75th Anniversary Volume)*, edited by G.F. Vander Voort et al., ASTM STP 1165, Philadelphia, 1993, pp. 266-273.
7. S.A. Saltykov, *Stereometrische metallographie*, Deutscher Verlag fur Grundstoffindustrie, Leipzig, 1974.
8. V.N. Vapnik, *Statistical learning theory*, John Wiley and Sons, New York, 1998.
9. C.C. Chang and C.J. Lin, LIBSVM: a library for support vector machines, 2001, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.