

## The influence of principal component analysis on the spatial structure of a multispectral dataset

C. RICOTTA, G. C. AVENA

Department of Plant Biology, University of Rome 'La Sapienza', Piazzale Aldo Moro 5, 00185 Rome, Italy

and F. VOLPE

ITALECO SpA, Via Carlo Pesenti 109, 00156 Rome, Italy

(Received 14 October 1997; in final form 28 September 1998)

**Abstract.** Multispectral remote sensing images often have extensive interband correlation. As a result, the images may contain similar information and have similar spatial structure. Principal component analysis (PCA) is a technique for removing or reducing the duplication or redundancy in multispectral images and for compressing all of the information that is contained in an original  $n$ -channel set of multispectral images into less than  $n$  channels or, more specifically, to their principal components. These are then used instead of the original data for image analysis and interpretation. The principal components are ranked in terms of the amount of variance that they explain. A consequence of ranking in this way is that the resulting principal components show a markedly different spatial structure from one another. This effect can be problematical, for example, when studying landscape ecology, where understanding the interactions between elements of the landscape structure as manifest in remote sensing images and environmental processes is of primary importance. Although the difference in spatial structure of an image after applying PCA and its influence on potential applications have been known for some time, it does not appear to have been studied explicitly. Accordingly, the aim of this paper was to examine the implications of applying PCA for the spatial structure and content of multispectral remote sensing images using parts of a Landsat Thematic Mapper (TM) frame of northern Sardinia, Italy. The results show that, due to the significant influence of PCA on the spatial structure and content of remote sensing images, the resulting principal components have a spatial structure and content that differ markedly from one another and from the original images. As a result, extreme care is necessary when applying PCA to remote sensing images and interpreting the results.

### 1. Introduction

In recent years, continuous advances in sensor technology have permitted the acquisition of remotely sensed data with improved spectral resolution. The detailed spectral resolution is helpful when mapping cover types, provided that the spatial resolution of the data is sufficient to represent a single surface type for each pixel (Price 1994). The selection of an appropriate spatial resolution is based both on the spatial structure of the analysed landscape and the desired scale of analysis (Marceau *et al.* 1994, Weishampel *et al.* 1994).

Terrestrial landscapes are not chaotic or random, but manifest spatial order. Landscape structure is composed of a mosaic of homogeneous patches arranged over the Earth's surface. Furthermore, landscape patches can be regarded as being organized in a nested hierarchy with each level containing the sub-levels below it (Forman and Godron 1986). For example, a forested landscape 'might be hierarchically composed of drainage basins, which in turn are composed of local ecosystems or stands, which in turn are composed of individual trees and tree gaps' (Forman 1995).

Remotely sensed images reflect in their properties the landscape patch hierarchy. The model proposed by Woodcock and Strahler (1987) assumes that remotely sensed images are composed of a nested hierarchy of discrete elements which are abstractions of landscape patches. Because of the complex nature of most environments, different classes of scene elements at a higher level in the nested hierarchy which are composed of completely different sets of lower-level components rarely exist. Conversely, higher-level scene elements might be composed of varying proportions of the same lower-level components. As this occurs, the internal variance associated with higher-level scene elements in the nested model will be high, reducing the spectral separability between different classes of higher-level scene elements and shifting the identification of elements on the remotely sensed image to a lower level in the nested scene model. In this sense, the influence of lower-level components on the internal variance of higher-level scene elements, often referred as 'scene noise', acts on remotely sensed images as any other unwanted source of noise, such as, for example, striping and bit errors (Woodcock and Strahler 1987).

Generally, the hierarchical level at which scene elements are identified on remotely sensed multispectral images is invariant to the spectral interval of observation. For instance, visual analysis of different bands/colour composites from a multispectral dataset with constant pixel resolution still reflects the same spatial structure, even if the contrast between different scene elements (i.e. forest patches versus non-forest patches) might considerably vary for the different band combinations (Bryant 1988). Due to the high similarity among individual bands of a multispectral image, statistical data compression tools like principal component analysis (PCA) are often applied in image analysis and image classification to reduce the amount of redundant information.

PCA undertakes a linear transformation of a set of numerical variables to create a new variable set with principal components that are uncorrelated and are ordered in terms of the amount of variance explained in the original data without *a priori* physical interpretation of the principal components (Eastman and Fulk 1993). Given a dataset with  $n$  variables,  $n$  principal components can be computed. With unstandardized PCA, each principal component is a linear combination of the original variables, with coefficients equal to the eigenvectors of the variance/covariance matrix. With standardized PCA, the eigenvectors are computed from the correlation matrix and the result is identical to standardizing all values (by subtracting the mean and dividing by the standard deviation) and computing the unstandardized principal components of the results. In both cases, the eigenvectors are customarily taken with unit length, and the sum of the variance in all the components is equal to the total variance present in the original input images.

In remote sensing applications, PCA has long been used as a data compression tool by discarding minor components with little explanatory value. For example, there is general agreement that the configuration of Landsat Thematic Mapper (TM)

data is well described in a three-dimensional Euclidean space, with minimum loss of total variance (Chavez and Howell 1988, Ceballos and Bottino 1997).

Regarding the spatial structure of remotely sensed data, a major consequence of PCA is that, unlike the original bands on which PCA is performed, the resulting principal components show a different spatial structure from one another. This effect can be particularly problematical in applications where understanding the interactions between the landscape structure as manifest on remote sensing images and the environmental processes is of primary importance. The objective of this study was to examine the implications of PCA on the spatial structure of remotely sensed multispectral images on the basis of the nested scene model proposed by Woodcock and Strahler (1987). Although the goal was essentially methodological, the following paragraphs provide a practical example based on the analysis of a Landsat 5 TM image of northern Sardinia, Italy.

## 2. Data

### 2.1. Study area

The area selected for investigation is the Olbian region on the northeast coast of Sardinia, Italy. This area is characterized by a typical Mediterranean climate, with cool winters and warm summers. The average annual rainfall is less than 500 mm, mainly concentrated during the autumn and spring, and it has a very dry summer.

The Olbian landscape has several kinds of patches arranged in a mosaic, the principal patches being fields combined with sparse cork oak stands and pasture lands. Natural vegetation is dominated by Mediterranean sclerophylls and includes all the transition forms of the Mediterranean vegetation, from the garrigue to the more evolved types of sclerophyllous woodlands. In the hilly areas, cork oak stands prevail, sometimes abandoned and undergoing transformation towards sclerophyllous woodlands.

### 2.2. Analytical procedures

A georeferenced Landsat TM 192/32 subscene of  $600 \times 600$  pixels acquired on 17 July 1995 was used for image analysis. The six standardized principal components (PC 1–6) were calculated from the transformation of the original TM reflective bands 1–5 and 7 (figures 1 and 2). The correlation matrix of the six original TM reflective bands and the percentage of variance that was mapped to each component after PCA are shown in tables 1 and 2, respectively.

To quantify the spatial structure of remotely sensed images, graphs of mean local variance as a function of spatial resolution can be used as proposed by Woodcock and Strahler (1987), to which readers are referred for further information. To place the current work in context, however, the characteristics of this measure will be outlined in order to clarify the relationships between spatial structure of remotely sensed images and the size of scene elements.

The data for these graphs is calculated by measuring the mean of the standard deviation computed from an odd-sized window of  $3 \times 3$  pixels and by this window moving across the image degraded to successively more coarse spatial resolutions. The reasoning behind this measure is as follows (Woodcock and Strahler 1987, Price 1992). If the scene elements are considerably larger than the pixel size, most of the pixels falling within the same scene element are highly correlated and local variance will be low. When the pixel resolution becomes similar to the dimension of the scene elements, the correlation among contiguous pixels decreases and local variance rises.

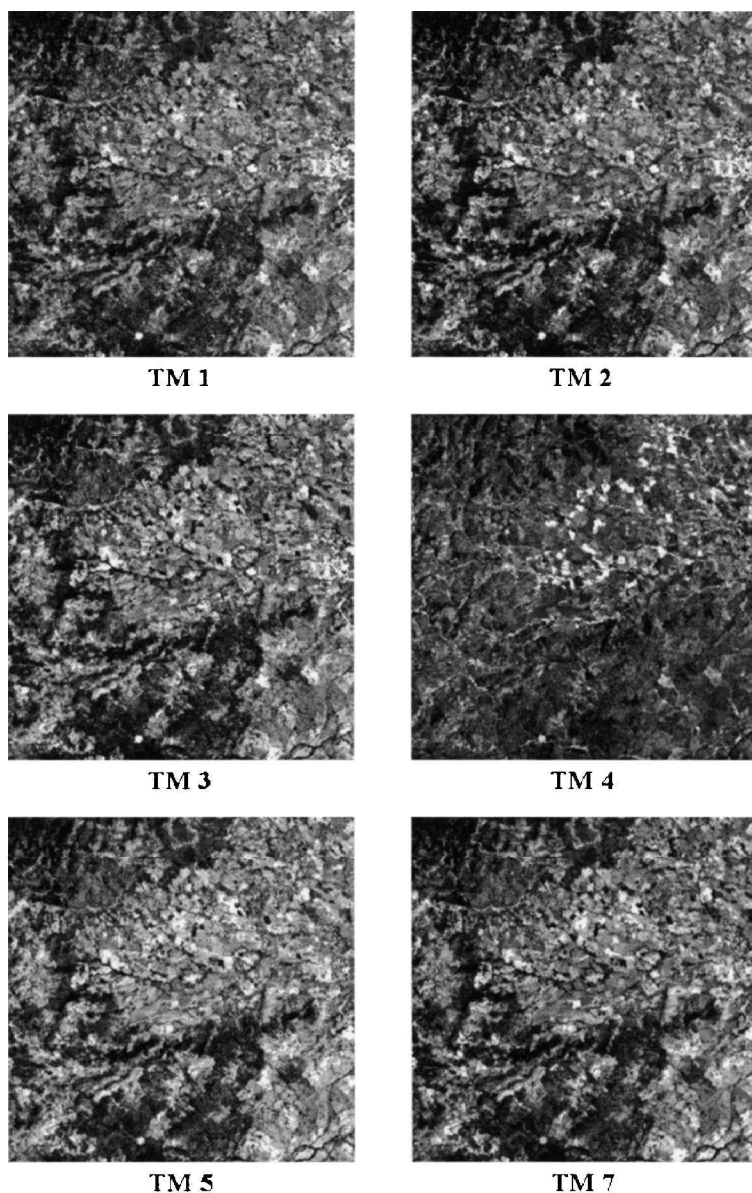


Figure 1. Landsat TM reflective bands of the study area. A 1 and 99% linear stretch was applied to all TM bands so that they would have approximately the same amount of contrast for the visual analysis. Note that, although there are differences, especially in the TM band 4, the TM reflective bands appear broadly similar in spatial structure.

At coarser resolutions, when the scene elements increasingly become smaller than the pixel size, all pixels look similar and local variance decreases again. To measure mean local variance at multiple spatial resolutions, the image data are degraded by combining original pixels into larger resolution cells and replacing the original digital numbers in each cell with their mean. Note that this degradation procedure of averaging original pixels over larger areas is incorrect in geostatistical terms (Cohen

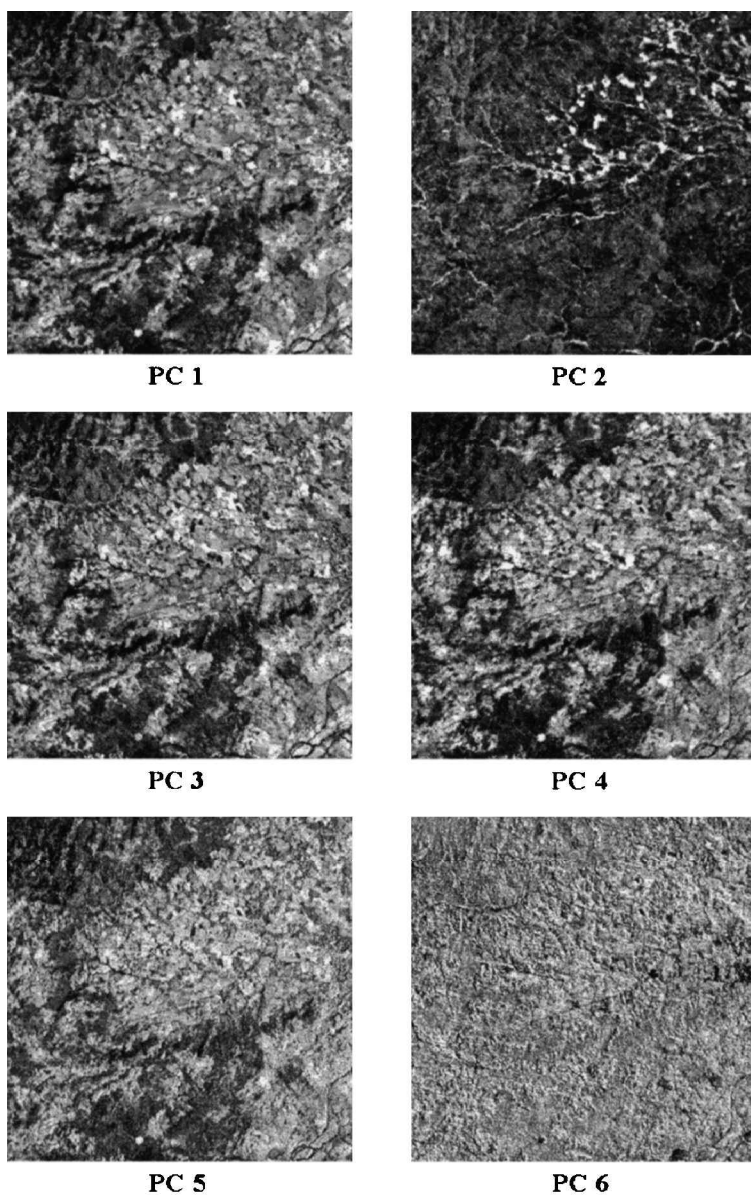


Figure 2. Principal components obtained from the original Landsat TM reflective bands of the study area. A 1 and 99% linear stretch was applied to all principal components so that they would have approximately the same amount of contrast for the visual analysis.

*et al.* 1990). However, the aim of this approach is not to replicate the response of a particular sensor, but to understand the effect of PCA on the spatial structure of remotely sensed images.

### 3. Results and discussion

The graphs of mean local variance as a function of spatial resolution for the six TM reflective bands of the study area are shown in figure 3. Note that, for each TM

Table 1. Correlation matrix of the Landsat TM reflective bands of the study area.

	TM 1	TM 2	TM 3	TM 4	TM 5	TM 7
TM 1	1.000	0.970	0.962	0.405	0.883	0.909
TM 2		1.000	0.977	0.468	0.870	0.899
TM 3			1.000	0.383	0.900	0.917
TM 4				1.000	0.382	0.322
TM 5					1.000	0.959
TM 7						1.000

Table 2. Percentage of variance mapped to each principal component of the study area.

Percentage of variance	
PC 1	81.61
PC 2	13.62
PC 3	3.29
PC 4	0.68
PC 5	0.59
PC 6	0.22

band, mean local variance at a given spatial resolution depends on the global variance of the TM band at original resolution. Thus, from a numerical point of view, mean local variance values of a single TM band can only be compared with other mean local variance values obtained from the same band degraded at different resolutions (Woodcock and Strahler 1987).

Despite this major limitation, the general shape of the graphs of figure 3 look similar as a result of the invariance of scene structure to the spectral range of observation. At the original resolution of the Landsat image, pixels are much smaller than scene elements. If a pixel is part of a distinct scene element, its neighbours are also likely to be in the same element and have similar values. In this situation, the mean local variance is generally low. As the resolution becomes more coarse, the likelihood that surrounding pixels will be similar decreases. Consequently, mean local variance increases. This trend continues until a peak is reached between 150 and 180 m, or 5–6 times the original resolution.

Minor differences occur between the graphs for TM bands 1–3 and 7, which peak at 150 m spatial resolution, whereas for TM bands 4 and 5 the peak is at 180 m. This small difference in the general shape of the graphs is principally due to the reduced contrast between vegetated and non-vegetated scene elements in the near-infrared bands with respect to the visible and mid-infrared bands (Chavez 1992). Because of the lack of extreme contrast between vegetated and non-vegetated elements in the near-infrared bands, the local detail often shows up better than in the visible and mid-infrared bands, causing a shift in the location of the peak in their graphs of mean local variance as a function of spatial resolution. After a peak is

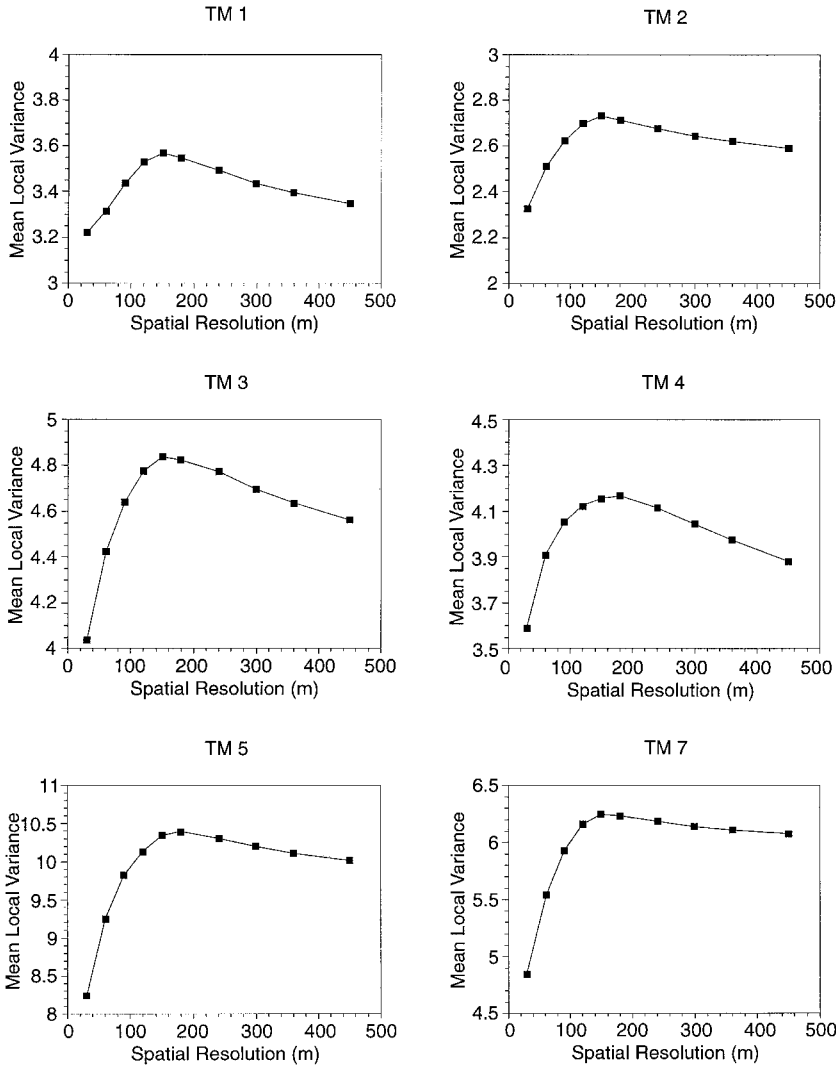


Figure 3. Graphs of mean local variance as a function of spatial resolution for the six Landsat TM reflective bands of the study area.

reached, mean local variance generally declines over the remaining resolutions as different elements are aggregated into coarser pixels.

The graphs of mean local variance for the six standardized principal components obtained from the transformation of the TM reflective bands are shown in figure 4. In this case, visual analysis shows that the resulting graphs reveal different patterns. The graph of mean local variance for PC 1 looks similar to the corresponding graphs of figure 3, showing a peak in mean local variance at 150m spatial resolution. In contrast, for PC 2, the peak in mean local variance shifts to 90 m, or just three times the original resolution. Furthermore, for PC 3–6, which are characterized by a residual variance varying from 3.29 to 0.22% (table 2), the graphs of mean local variance start with high mean local variance at initial resolution and decline rapidly as pixel size increases.

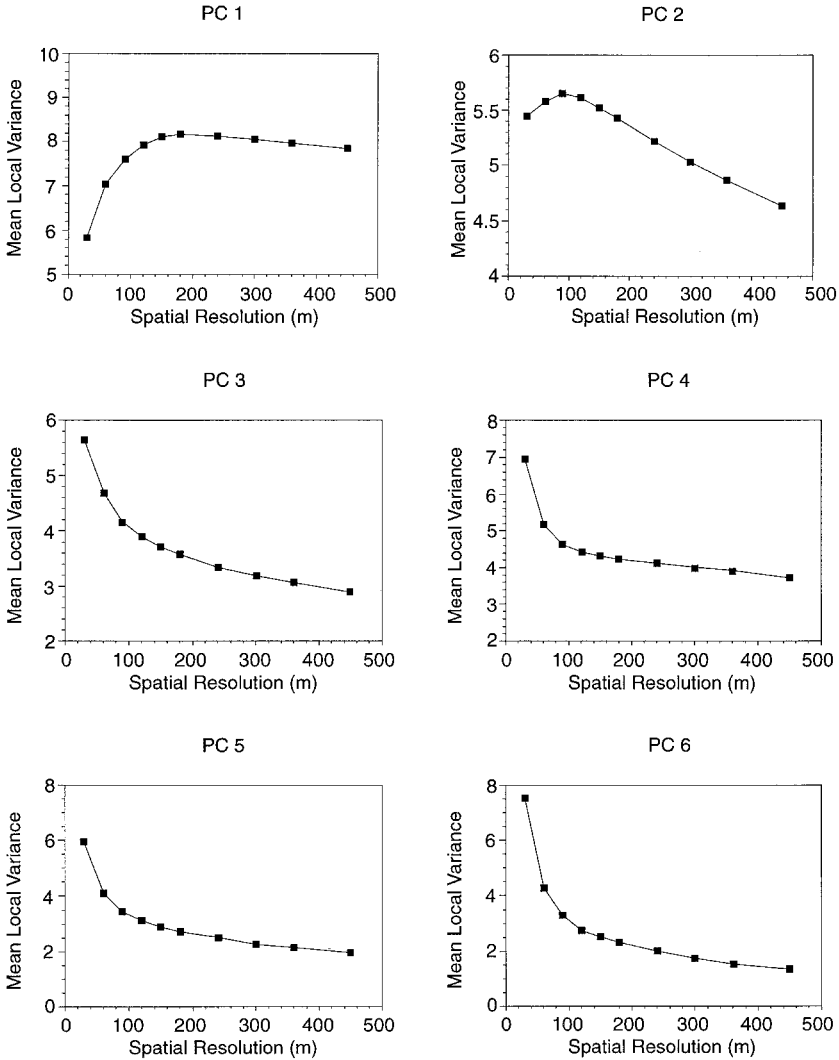


Figure 4. Graphs of mean local variance as a function of spatial resolution for the six principal components obtained from the original Landsat TM reflective bands of the study area.

To understand the perspective that different principal components show different spatial structures from one another, an understanding of the effects of PCA on the hierarchical level at which scene elements are identified on remotely sensed multispectral images is required. The major consequence of ranking principal components in terms of the amount of variance explained is that the influence of noise becomes increasingly important in the higher-order principal components which are dominated by residual variance (Green *et al.* 1988). For remotely sensed images, the chief source of variance is caused by differences among different classes of higher-level elements in the nested scene model (Woodcock and Strahler 1987). Therefore, the spatial structure of PC 1, which accounts for 81.61% of total variance (table 2), is similar to the corresponding structure of the original TM reflective bands. Conversely,



the major effect of noise on the spatial structure of the higher-order principal components increases the internal variance of higher-level scene elements, shifting the identification of scene elements to increasingly lower hierarchical levels in the nested model. Note that, for PC 3–6, mean local variance reaches its maximum at the original resolution of the TM reflective bands, giving way to a one-pixel-one-element spatial structure.

#### 4. Conclusions

This paper examined the implications of PCA on the spatial structure of a Landsat TM multispectral image of northern Sardinia, Italy, on the basis of the nested scene model in remote sensing proposed by Woodcock and Strahler (1987).

The consequence of ranking principal components in terms of the amount of variance explained is that, unlike the original bands on which PCA is performed, the resulting principal components show a different spatial structure from one another. This effect can be particularly problematical in topics like landscape ecology where understanding the interactions between the remotely sensed landscape structure and the environmental processes is of primary importance. Furthermore, besides PCA, the same effect might be evident in all  $n$ -space indices (Perry and Lautenschlager 1984, Elvidge 1985) which are based on the redistribution of image variance of multispectral datasets, such as, for example, the Tasseled Cap (TC) transformation (Crist and Cicone 1984). For instance, TC is similar to PCA. As such, a TM image can be reduced from a set of six reflective bands to the first three TC axes, namely Brightness, Greenness and Wetness, without much loss of information (Cohen 1994). However, by choosing the order in which factorization of the total variance in the dataset is performed, the TC transformation imposes a physical interpretation on the resulting axes.

#### Acknowledgments

The authors wish to thank the anonymous referees for the stimulating comments on a previous version of this paper.

#### References

- BRYANT, J., 1988, On displaying multispectral imagery. *Photogrammetric Engineering and Remote Sensing*, **54**, 1739–1743.
- CEBALLOS, J. C., and BOTTINO, M. J., 1997, The discrimination of scenes by principal component analysis of multi-spectral imagery. *International Journal of Remote Sensing*, **18**, 2437–2449.
- CHAVEZ, P. S., 1992, Comparison of spatial variability in visible and near-infrared spectral images. *Photogrammetric Engineering and Remote Sensing*, **58**, 957–964.
- CHAVEZ, P. S., and BOWELL, J. A., 1988, Comparison of the spectral information content of Landsat Thematic Mapper and SPOT for three different sites in the Phoenix, Arizona Region. *Photogrammetric Engineering and Remote Sensing*, **54**, 1699–1708.
- COHEN, W. B., 1994, GIS application perspective: current research on remote sensing of forest structure. In *Remote Sensing and GIS in Ecosystem Management*, edited by V. A. Sample (Washington, DC: Island Press), pp. 91–107.
- COHEN, W. B., SPIES, T. A., and BRADSHAW, G. A., 1990, Semivariograms of digital imagery for analysis of conifer canopy structure. *Remote Sensing of Environment*, **34**, 167–178.
- CRIST, E. P., and CICONE, R. C., 1984, A physically-based transformation of Thematic Mapper data—the TM Tasseled Cap. *IEEE Transactions on Geoscience and Remote Sensing*, **22**, 256–263.
- EASTMAN, J. R., and FULK, M., 1993, Long sequence time series evaluation using standardized principal components. *Photogrammetric Engineering and Remote Sensing*, **59**, 991–996.

- ELVIDGE, C., 1985, Influence of rock-soil spectral variation on the assessment of green biomass. *Remote Sensing of Environment*, **17**, 265–279.
- FORMAN, R. T. T., 1995, *Land Mosaics* (Cambridge: Cambridge University Press).
- FORMAN, R. T. T., and GODRON, M., 1986, *Landscape Ecology* (New York: John Wiley).
- GREEN, A. A., BERMAN, M., SWITZER, P., and MAURICE, D. C., 1988, A transformation for ordering multispectral data in terms of image quality with implications for noise removal. *IEEE Transactions on Geoscience and Remote Sensing*, **26**, 65–74.
- MARCEAU, D. J., GRATTON, D. J., FOURNIER, R. A., and FORTIN, J. P., 1994, Remote sensing and the measurement of geographical entities in a forested environment. 2. The optimal spatial resolution. *Remote Sensing of Environment*, **49**, 105–117.
- PERRY, C. R., and LAUTENSCHLAGER, L. F., 1984, Functional equivalence of spectral vegetation indices. *Remote Sensing of Environment*, **14**, 169–182.
- PRICE, J. C., 1992, Estimating vegetation amount from visible and near infrared reflectances. *Remote Sensing of Environment*, **41**, 29–34.
- PRICE, J. C., 1994, How unique are spectral signatures? *Remote Sensing of Environment*, **49**, 181–186.
- WEISHAMPEL, J. F., SUN, G., RANSON, K. J., LEJEUNE, K. D., and SHUGART, H. H., 1994, Forest textural properties from simulated microwave backscatter: the influence of spatial resolution. *Remote Sensing of Environment*, **47**, 120–131.
- WOODCOCK, C. E., and STRAHLER, A. H., 1987, The factor of scale in remote sensing. *Remote Sensing of Environment*, **21**, 311–332.