# Hyperspectral Image Database Query Based on Big Data Analysis Technology

Qi Beixun[1]

[1]Engineering from Jilin University，Jilin,China,130022.
(École centrale de Nantes) France；

**Abstract:** In this paper, we extract spectral image features from a hyperspectral image database, and use big data technology to classify spectra hierarchically, to achieve the purpose of efficient database matching. In this paper, the LDMGI (local discriminant models and global integration) algorithm and big data branch definition algorithm are used to classify the features of the hyperspectral image and save the extracted feature data. Hierarchical color similarity is used to match the spectrum. By clustering colors, spectral information can be stored as chain nodes in the database, which can improve the efficiency of hyperspectral image database queries. The experimental results show that the hyperspectral images of color hyperspectral images are highly consistent and indistinguishable, and need to be processed by the machine learning algorithm. Different pretreatment methods have little influence on the identification accuracy of the LDMGI model, and the combined pretreatment has better identification accuracy. The average classification accuracy of the LDMGI model training set is 95.62%, the average classification accuracy of cross-validation is 94.36%, and the average classification accuracy of the test set is 89.62%. Therefore, using big data analysis technology to process spectral features in hyperspectral image databases can improve query efficiency and more accurate query results.

## 1 Introduction

In recent years, deep learning has made significant progress in the hyperspectral image classification task, but manual labeling of hyperspectral images is still time-consuming and laborious, and deep learning methods usually rely on a large number of manually labeled samples, when only a small number of labeled samples, the deep learning model is difficult to train. Therefore, using the deep learning method for hyperspectral image classification is still faced with the problem of lack of labeled training samples and deep neural network training.

At present, many scholars have researched image data sets. For example, Xie t has collected a set of data sets called palette and text, which contains more than 10000 pairs of color themes and their annotation texts [1]. Then, Lee s uses the data set to train a set of conditions to generate the countermeasure network and converts the database records input by users into several groups of color theme outputs containing five colors [2]. For example, Zhang l proposed a color theme extraction method based on a regression model [3]. Recently, Liy proposed a geometric-based color theme extraction method, which uses convex hull in RGB color space to represent a color theme, and transforms the problem of color theme extraction into the problem of convex hull

generation and simplification in geometric space [4]. On this basis, fan h further proposes an improved method based on iterative optimization, which improves the problem of poor color representation in color themes [5].

Different from the above methods, Xing l proposes a color theme extraction method for multiple images, which can maintain the color consistency between different images [6]. The extracted color theme can be used in many different applications, such as image color enhancement and image re-coloring. In addition, Baloch g proposed a machine learning-based color theme aesthetic evaluation method [7]. Dabov K uses 24 pre-defined color emotions to achieve image emotion transfer [8]. By analyzing a data set containing 1000 color combinations, Othman h realizes image index based on color emotion [9]. Luc implements image color editing based on the database record concept, but only supports 15 predefined database record concepts [10]. However, the above research was done spontaneously by volunteers on the Internet, and the labeling task has certain requirements on the color matching level and aesthetics of the taggers. It is difficult to guarantee the quality of training data under the condition that the self standard of the tagger cannot be ensured.

In this paper, we extract spectral image features from a hyperspectral image database and use big data technology to classify spectra hierarchically, so as to achieve the purpose of efficient database matching. In this paper, the LDMGI algorithm and big data branch

definition algorithm are used to classify the features of the hyperspectral image and save the extracted feature data. Hierarchical color similarity is used to match the spectrum. By clustering colors, spectral information can be stored as chain nodes in the database, which can improve the efficiency of a hyperspectral image database query.

## 2 Hyperspectral Adaptive Database and Image Labeling Method

### 2.1 Hyperspectral Image Labeling Method Based on Deep Learning

Using the deep learning method to classify hyperspectral images is still faced with the problem of lack of labeled training samples and deep neural network training [11]. A large number of studies and practices have shown that the features extracted from spatial-spectral features and depth models can help to improve the classification accuracy of hyperspectral images [12]. For this reason, in order to solve the problem that the network is difficult to train and the accuracy is not high when there are only a few labeled samples in deep learning, the author designs a deep forest classification method of the hyperspectral image combined with spatial-spectral information [13]. The designed method first extracts spatial-spectral features to improve the classification accuracy of hyperspectral images, and then in order to further improve the classification accuracy and avoid the problem that the deep learning model is difficult to train when the number of labeled samples is small, the spatial-spectral features are input into the deep forest model for classification [14]. The deep forest model can automatically extract deep features and has the advantages of fewer parameters and easy training. In this way, the deep features are effectively utilized, and the problem that the deep learning model is difficult to train is avoided [15].

### 2.2 Hyperspectral Adaptive Database

The adaptive representation of hyperspectral data is used to construct a graph to reveal the hyperspectral characteristics of data. Multi structure manifold embedding designs hyperspectral maps and hypermaps through hyperspectral representation and reveals the inherent hyperspectral structure in high-dimensional data by combining these maps. However, the hyperspectral representation model based on L1 normal form produces unstable hyperspectral representation due to the potential instability of hyperspectral decomposition. In order to alleviate this instability, some improved hyperspectral representation models are applied to hyperspectral image feature extraction. A constrained hyperspectral image is smoothed by adding a manifold regularization term to hyperspectral optimization. A new hyperspectral optimization model is designed based on local manifold discriminant learning. The similarity weight based on euclidean distance is used to constrain the hyperspectral

representation model to reveal the hyperspectral structure of data based on the local manifold. Although the above hyperspectral representation model makes the hyperspectral decomposition more stable by imposing manifold constraints, the manifold constraints depend on the quality of the predefined similarity measure matrix. The similarity measure matrix predefined by Gaussian kernel function is usually affected by noise and other factors, so it is difficult to design an appropriate predefined similarity measure matrix. To solve these problems, this paper proposes a local discriminant and global hyperspectral preserving projection algorithm based on Hyperspectral representation and learning graph regularization. Based on the manifold regularization term, a hyperspectral representation model with learning graph regularization constraints is proposed, and the locally linear structure is obtained adaptively, which solves the problem of LDMGI (local discriminant models and global integration). When the k-nearest neighbor graph is used in the algorithm, it is difficult to select neighborhood parameters. At the same time, the algorithm can also use the hyperspectral representation model based on learning graph regularization to extract the global hyperspectral structure in the data. By combining local discriminant information with global hyperspectral structure, discriminant features including both local and global features are extracted, so as to improve the final classification performance. This method uses a set of determining strategies to select the initial value of the K-means clustering method. First, select the color corresponding to bucket J with the largest number of pixels JN. Then, the pixel number in of each bucket I is attenuated proportionally, and the corresponding color of the bucket with the largest number is selected

$$f(\mathrm{JN\text{-}H}) = \frac{1}{\mathrm{J}\text{-}h}\sum_{i=1}^{N} k(\frac{X_i - (\mathrm{I\text{-}H})}{\mathrm{IN}}) \qquad (1)$$

Where h is the distance from barrel I to barrel J in lab space; Jn is a parameter related to attenuation radius, which is set to K by default in the experiment. Repeat the above steps K times to get the initial value of color theme containing K colors

$$K(U) = \sigma t = \frac{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(I_{it} - I_{it})^2}}{I_{it}} \qquad (2)$$

$$X_t = \tanh(w_c \mathrm{K}_t + u_c(\mathrm{JN}_t \Theta h_{t-1}) + \mathrm{K}_c) \qquad (3)$$

In addition, there are large areas of black or white color patches in most real images, and the color is usually independent of the subject of the image. To solve the above problem, this paper adopts a set of flexible schemes, instead of the implementation of the clustering method, and fixed two of these colors as pure black and pure white, and removed them before the final output.

After the above steps, this method can extract n color themes from n images, and make each color theme contain K colors. As the data comes from the Internet image database or online search platform, most of the collected images are real-life images rather than synthetic images. Therefore, among these color themes, the color

theme mode with higher frequency may be expected by users. Because the colors in the color theme are unordered, two-color themes m and N need to be defined first

$$M_{nb} = \sum_{j=2}^{k} \sum_{h=1}^{j-1} G_{jh}(p_j s_h + p_h s_j) D_{jh} \qquad (4)$$

$$N = \frac{\sum_{j=1}^{k} \sum_{h=1}^{k} \sum_{t=1}^{n_j} \sum_{r=1}^{n_h} |y_{ij} - y_{hr}|}{2n^2 u} \qquad (5)$$

It is used to quantitatively describe the degree to which they belong to the same color theme. When defining this metric, we need to ensure that it satisfies not only the basic properties of the metric, but also the properties that the measurement results are independent of the order of colors in the color theme; It defines the metric of two color themes in the color theme space as

$$M/N = \sum_{T} = color(\max(\sigma_i - \upsilon, 0)) \qquad (6)$$

Where color represents the measurement of two colors in the color space. In addition, since the output result is only a single color. The range and diversity of database records supported by this method are better than this method. To sum up, the results of the qualitative analysis show that this method is slightly better than this method in terms of the correctness and diversity of output results and the acceptable input range.

## 3  Image Database Query Design

### 3.1 Methods

In this paper, we extract spectral image features from a hyperspectral image database and use big data technology to classify spectra hierarchically, to achieve the purpose of efficient database matching. In this paper, LDMGI algorithm and big data branch definition algorithm are used to classify the features of the hyperspectral image and save the extracted feature data. Hierarchical color similarity is used to match the spectrum. By clustering colors, spectral information can be stored as chain nodes in the database, which can improve the efficiency of hyperspectral image database queries.

### 3.2 Design

This paper proposes a method to automatically generate color themes from database records, which supports the generation of multiple different color themes from the same database record. In this paper, we propose to use the real image as a new data source to solve the problem of color theme generation based on database records. Thanks to more diversified data, this method can support more abundant database record input than existing methods, and even for some rare input, this method can calculate reasonable results. This paper puts forward a set of quality indicators so that the final output results can be arranged in the order of good to bad, thus saving the user's search time. This method involves the training of neural networks, and the number of super parameters generated in the process reduces the workload, so it is easier to debug and analyze, and more robust. This method can also provide more abundant training data sources for the existing deep learning-based methods and can be easily integrated with the neural network-based methods, ensuring the accuracy and diversity of the output results while ensuring the end-to-end characteristics and real-time.

## 4  Results and Discussion

This method supports user input in the form of any word or phrase. It is impossible to list all possible images as input and test them. To compare comprehensively and fairly, six different input images are randomly selected under each type, and the training results are shown in Figure 1.

Although the metric between any two color themes can be easily calculated, the normed space from which the metric can be derived cannot be obtained directly. This means that we can't directly perform linear operations on color themes, such as calculating the average value of a group of color themes, so we can't directly use k-means and other methods to achieve clustering. The average classification accuracy of the training set is 98%, the average classification accuracy of the cross-validation is 97.52%, and the average classification accuracy of the test set is 92.32%. In addition, for this step, it is difficult to directly determine the number of categories to be clustered, because it is closely related to the distribution of color themes in n images. In this paper, the nearest neighbor propagation method is used to solve the clustering problem. N color topics are regarded as N nodes in a graph, and only the similarity between them is considered, without any assumption of their specific algebraic structure. By alternately calculating and transferring the attraction degree and attribution degree between nodes, the method can quickly converge and get the color theme contained in each group.
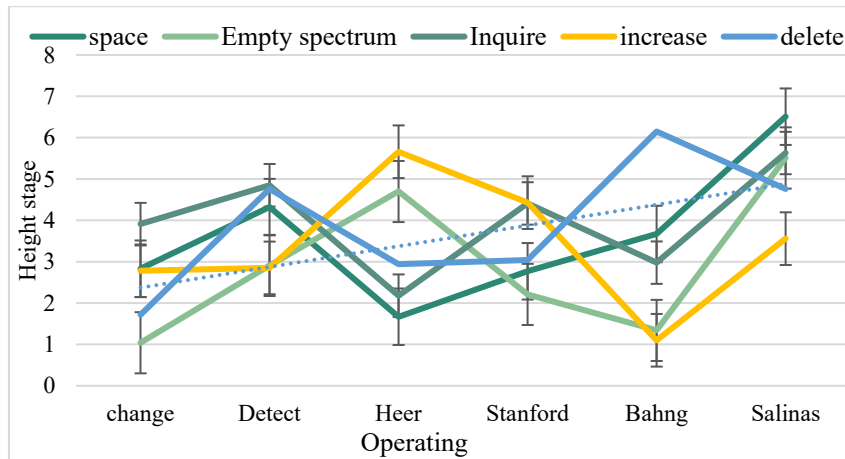
**Figure 1.** Input image of hyperspectral image

**Table 1.** Attraction and attribution between nodes

| Item | Inquire | increase | delete | change |
|------|---------|----------|--------|--------|
| Heer | 1.31 | 1.62 | 1.3 | 1.7 |
| Stanford | 4.07 | 4.26 | 5.99 | 5.3 |
| Bahng | 2.99 | 3.07 | 3.49 | 5.46 |
| Salinas | 4.74 | 2.35 | 3.52 | 3.44 |
| spectrum | 6.95 | 5.19 | 6.52 | 2.05 |

As shown in Table 1, the hyperspectral images of color hyperspectral images are highly consistent in shape, and the distinction is not obvious, so it needs to be processed with the help of a machine learning algorithm. Different pretreatment methods have little influence on the identification accuracy of LDMGI model, and the combined pretreatment has better identification accuracy. The average classification accuracy of LDMGI model training set is 95.62%, the average classification accuracy of cross-validation is 94.36%, and the average classification accuracy of the test set is 89.62%. The average classification accuracy of the training set is 98%, the average classification accuracy of the cross-validation is 97.52%, and the average classification accuracy of the test set is 92.32%.
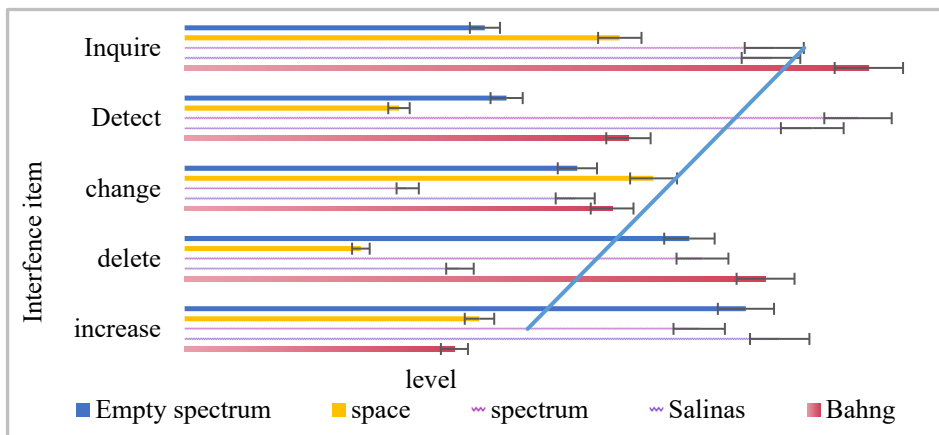


**Figure 2.** Image information contained in database records

As shown in Figure 2, the two groups of results were analyzed qualitatively and quantitatively in terms of correctness and diversity. Among them, correctness measures whether the generated color themes are close enough to the image information contained in the input database records, while diversity measures whether the generated color themes are rich enough, that is, whether the range of color themes is wide enough. In addition, in order to improve the comprehensiveness of comparative analysis, qualitative comparative experiments are carried out between this method and Heer's method. However, due to the limited input of database records supported by Heer and other methods, the number of examples that can be processed normally is small, and the quantitative results are not statistically significant.
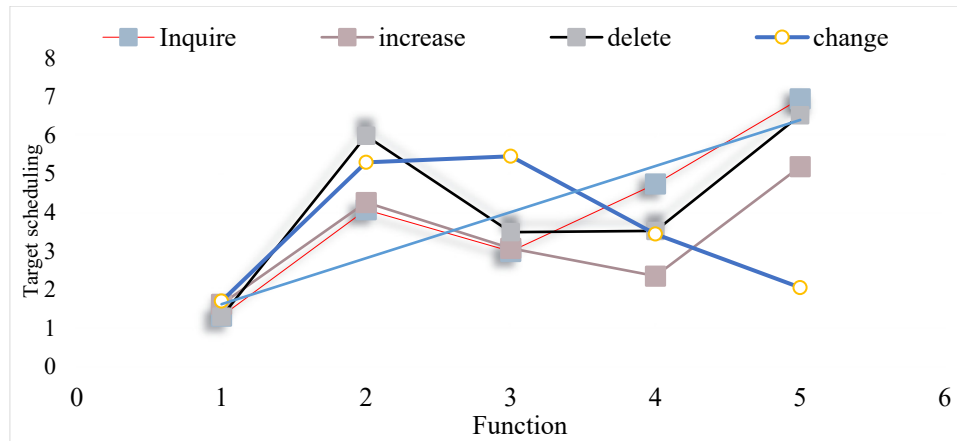
**Figure 3.** Partially input database records

Due to the lack of consideration of the internal image information and structure of the image data, the performance of this method is still poor for some input database records, as shown in Figure 3. On the one hand, most of the area in the input image corresponds to the sky or land, so without saliency detection, the color theme extracted directly must be the sky and land, not the spectrum itself. On the other hand, because the spectrum is translucent, to extract the color theme of the spectrum accurately, further soft segmentation is needed to separate or remove the color of the background sky or land. Therefore, to further improve the quality of output results, the introduction of saliency detection and image soft segmentation is indispensable. This will be the future work of this paper. In terms of time efficiency, there is still room for improvement.

## 5  Conclusions

In this paper, LDMGI algorithm and big data branch definition algorithm are used to classify the features of the hyperspectral image and save the extracted feature data. Hierarchical color similarity is used to match the spectrum. By clustering colors, spectral information can be stored as chain nodes in the database, which can improve the efficiency of hyperspectral image database queries. This method can get relatively acceptable results. Among the first five-color themes, three-color themes contain red color. In addition, because the images supported by other methods are very limited, the test input images involved in the experiment can not be processed, so this paper manually replaces the input image with the image input which is the closest in semantics and can be accepted by this method. Therefore, using big data analysis technology to process spectral features in hyperspectral image databases can improve query efficiency and more accurate query results.

## Reference

1. XIE T, LI S, SUN B. Hyperspectral images denoising via nonconvex regularized low-rank and sparse matrix decomposition[J]. IEEE Transactions on Image Processing, 2019, 29(1): 44-56.

2. LEE S, NEGISHI M, URAKUBO H, et al. Mu-net: Multi-scale U-net for two-photon microscopy image denoising and restoration[J]. Neural Networks, 2020, 125(5): 92-103.

3. ZHANG L, WANG J, AN Z. Classification method of CO2 hyperspectral remote sensing data based on neural network[J]. Computer Communications, 2020,156(5): 124-130.

4. LI Y, XU J, XIA R, et al. A two-stage framework of target detection in high-resolution hyperspectral images[J]. Signal, Image and Video Processing, 2019, 13(7): 1339-1346.

5. FAN H, LI J, YUAN Q, et al. Hyperspectral image denoising with bilinear low rank matrix factorization[J]. Signal Processing, 2019, 163: 132-152.

6. XING L, CHANG Q, QIAO T. The algorithms about fast non-localmeans based image denoising[J]. Acta Mathematicae Applicatae Sinica, English Series, 2019, 28(2): 247-254.

7. BALOCH G, OZKARAMANLI H. Image denoising via correlation-based sparse representation[J]. Signal, Image and Video Processing, 2017, 11(8): 1501-1508.

8. DABOV K, FOI A, KATKOVNIK V, et al. Image denoising by sparse 3-D transform-domain collaborative filtering[J]. IEEE Transactions on image processing, 2017, 16(8): 2080-2095.

9. OTHMAN H, QIAN S E. Noise reduction of hyperspectral imagery using hybrid spatial-spectral derivative-domain wavelet shrinkage[J]. IEEE Transactions on Geoscience and Remote Sensing, 2019, 44(2): 397-408.

10. LU C, TANG J, YAN S, et al. Nonconvex nonsmooth low rank minimization via iteratively reweighted nuclear norm[J]. IEEE Transactions on Image Processing, 2019, 25(2): 829-839.

11. LI C, MA Y, HUANG J, et al. Hyperspectral image denoising using the robust low-rank tensor recovery[J]. Journal of the Optical Society of America A, 2019, 32(9): 1604-1612.

12. RENARD N, BOURENNANE S, BLANC-TALON

J. Denoising and dimensionality reduction using multilinear tools for hyperspectral images[J]. IEEE Geoscience and Remote Sensing Letters, 2018, 5(2): 138-142.

13. KONG X, ZHAO Y, XUE J, et al. Hyperspectral Image Denoising Using Global Weighted Tensor Norm Minimum and Nonlocal Low-Rank Approximation[J]. Remote Sensing, 2019, 11(19): 2281-2303.

14. ZHANG H, HE W, ZHANG L, et al. Hyperspectral image restoration using low-rank matrix recovery[J]. IEEE transactions on geoscience and remote sensing, 2019, 52(8): 4729-4743.

15. CANDÈS E J, LI X, MA Y, et al. Robust principal component analysis?[J]. Journal of the ACM, 2019, 58(3): 31-37.