

Image Indexing and Retrieval using Cross-Entropy Measures

M. Srikanth

Department of Computer Science and Engineering
State University of New York at Buffalo, Amherst, NY, USA. 14228-2567
srikanth@cedar.buffalo.edu

Abstract— In content-based image indexing and retrieval, images are indexed based on the image attributes like, color, texture, shape. These attributes are used in enabling image retrieval based on query images. A user provides a sample image and the images in the database that are similar to the sample image are retrieved. Image content is represented in the form of feature vectors or histograms. In this paper, color content of the image represented in the form of probability density functions is used as image signatures for image indexing and retrieval. Similarity measure based on cross-entropy measure, specifically Kullback-Leibler measure, defined on probability density functions is used for comparing the signature of the images. Some results of image indexing and retrieval are presented.

I. INTRODUCTION

The performance of an image retrieval system depends on the image attributes used in indexing an image database as well as the representation of the indices. Current image retrieval techniques use attributes like, color, texture, shape of objects to characterize images. These attributes are represented as feature vectors or histograms. While texture and shape are stored as feature vectors, color content is usually characterized by a histogram of its distribution.

In feature vector representation of image attributes, each component of the feature vector represents a *single measurement* taken over the entire image(global features) or a region of interest(local features). A distance measure is used to compare feature vectors, thereby allowing the retrieval of images with “similar” image attributes. Histograms, on the other hand, represent the *distribution of a particular feature* in the image. Although the feature vectors provide a simpler representation of the image content, it usually represents the “average” characteristic of the image. The histogram, represented by higher dimensional data, provides more information about the image.

As a trade off between efficiency and details, probability density functions(PDFs) have been explored to represent image content[6]. Kelly et al.[5] present results of their tests on using PDF representation for color content on the Red-Green-Blue(RGB) color space in the indexing and retrieval of images from Smithsonian Image Database. They use the Euclidean distance between the probability density

functions as a similarity measure for comparing the images.

This paper reports on the results of using probability density functions for representing color content of images in the *Hue-Value-Chroma*(HVC) color space and using the Kullback-Leibler cross-entropy measure as the similarity measure for comparing the images. The HVC color space forms the basis for indexing the images in a database. The HVC features of an image can be obtained by a non-linear transformation of its RGB values. The HVC color space has been shown to be effective in detecting regions of a given color[3] in an image and is insensitive to lighting variations in the image.

Given a query image, its signature in the form of probability density function is computed and compared with the signature of the images in the database using the Kullback-Leibler cross entropy measure. The result of image retrieval is an ordered list of the images in the database based on the cross-entropy distance. Some results of image retrieval from the KODAK image database are presented in section IV.

II. IMAGE INDEXING

The indexing of images in an image database is the pre-processing step of preparing the image database to support content-based image retrieval. This involves generating an index or a “signature” for each image in the database. Two significant steps in indexing an image database are (1) identifying and extracting image features for indexing and (2) identifying a representation of the extracted features as indices for the images.

The color content of the image represented in the HVC color space is used here as the image attribute for indexing the database. A probability density function is used to represent the HVC features of an image. It is noted that unlike other image retrieval techniques that use discrete valued histograms, this approach uses continuous probability density functions to characterize the distribution of colors in an image. Although it is more expensive and difficult to compute continuous probability density functions than histograms, it provides a measure of “closeness” of two specific colors(or features) to each other. In histograms, each color is assigned to a specific bin. All colors assigned to a particular bin are considered to be equivalent and dissimilar from colors assigned to a different bin.

This work was funded in part by Eastman Kodak Company.

Published in the Proceedings of the HKK Conference, Waterloo, Ontario, Canada, June 1999.

A. HVC color space

The Hue, Value and Chroma components of an image are computed at every pixel in the image. The red-green-blue (RGB) color model is the most widely used model to represent images because of its ease of implementation. However, the HVC space is most suited for identifying regions of a particular color in an image. The transformation from RGB color space to HVC color space is non-linear.

In the RGB color space, each pixel is represented in 24-bits — 8-bits each for R, G, and B, respectively. The following steps are involved in the conversion of RGB to HVC values [3].

$$X = 0.607R + 0.174G + 0.201B \quad (1)$$

$$Y = 0.299R + 0.587G + 0.114B \quad (2)$$

$$Z = 0.066G + 1.117B \quad (3)$$

$$L^* = 116 \left(\frac{Y}{Y_0} \right)^{1/3} - 16 \quad (4)$$

$$a^* = 500 \left[\left(\frac{X}{X_0} \right)^{1/3} - \left(\frac{Y}{Y_0} \right)^{1/3} \right] \quad (5)$$

$$b^* = 200 \left[\left(\frac{Y}{Y_0} \right)^{1/3} - \left(\frac{Z}{Z_0} \right)^{1/3} \right], \quad (6)$$

where X_0, Y_0, Z_0 are $X, Y,$ and Z values for the reference white color, respectively.

The HVC values are derived from $L^*a^*b^*$ values with the equations

$$H = \arctan(b^*/a^*) \quad (7)$$

$$V = L^* \quad (8)$$

$$C = \sqrt{(a^*)^2 + (b^*)^2}. \quad (9)$$

B. PDF representation

The distribution of the HVC features in the image is represented by a probability density function. This is the content signature of the image. There are several techniques for probability density function estimation in the literature[7]. Here, the probability density function is modeled as a linear combination of Gaussian distributions and the images in the database are used to define the parameters for the Gaussian mixture.

A mean vector \mathbf{m}_i and covariance matrix R_i characterize a Gaussian distribution. A general data clustering of the feature vector of the images can be used to generate clusters for which the \mathbf{m}_i and R_i can be obtained. Any cluster having a singular covariance matrix is deleted and ignored in subsequent processing since the computation of distance measure involves inverting the covariance matrix.

The next step in generating content signatures for images is computing the weights for the Gaussian distributions corresponding to the clusters. The Gaussian distributions are weighted by the number of pixels in the image belonging to

the corresponding cluster. Let K be the number of clusters and let N be the number of images in the image database. The signature of the i th image in the image database is given by

$$P_i(\mathbf{x}) = \sum_{j=1}^K w_{ij} G_j(\mathbf{x}) \quad (10)$$

where

$$G_j(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |R_j|^{-1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{m}_j)^T R_j^{-1} (\mathbf{x} - \mathbf{m}_j) \right]. \quad (11)$$

The collection of Gaussians $G = \{G_1, G_2, \dots, G_K\}$, defined by their mean vectors and covariances matrices, is fixed for an image database. The weight vector for the i th image is given by $W_i = \{w_{i1}, w_{i2}, \dots, w_{iK}\}$. The pair, $\{G, W_i\}$, of Gaussian mixture and weight vector is the signature of the image.

For clarity, the following notation is used to differentiate the signatures of two different images, I_1 and I_2 :

$$P_{I_1}(\mathbf{x}) = \sum_{j=1}^K w_j G_j(\mathbf{x}) \quad P_{I_2}(\mathbf{x}) = \sum_{j=1}^K v_j G_j(\mathbf{x}) \quad (12)$$

The image indexing is performed in two phases. In the first phase the cluster centers for the Gaussian mixture is estimated and in the next phase the signatures for all the images in the database are computed.

C. Generate Cluster Centers

The cluster centers are the mean vectors of the Gaussian distributions in the mixture. It is assumed that there are K distributions in the mixture. The feature vectors in the HVC color space are obtained using the RGB to HVC transformation given in section II-A. A general data clustering algorithm can be used for determining the cluster centers. The K -means clustering and unsupervised Bayesian learning methods were considered [1]. It is reasonable to assume that an image database is not a fixed collection of images. Although the unsupervised Bayesian learning method is well suited for this problem, the K -means clustering method, with a small modification, is used for computing the cluster centers.

In K -means clustering, the HVC values at each pixel for all the images in the image database have to be computed and made available during the computation of the cluster centers. This requires significant computing resources, space and time. As an alternative, a subset of images in the database are used to compute the cluster centers. The clustering thus obtained is only an approximation to the clustering of the data from all the images in the database.

D. Generating Image Signatures

The signature of an image is given by the Gaussian mixture representing the clusters in the image database and a weight vector corresponding to the cluster distribution in

the image. The Gaussian mixture is represented by a set of Gaussian distributions characterized by their mean vectors and covariance matrices.

Using the cluster centers generated in the first stage, a signature is generated for each image in the image database. For each pixel in an image, the feature vector is extracted, it is assigned to cluster C_i , where the cluster center of C_i is the closest point to the feature vector using the L_2 norm. This classification is used to update the covariance matrix of the clusters.

The number of pixels in an image that are assigned to a particular cluster is the corresponding weight of the cluster in the image. The signature of an image is a vector of weights attached to the clusters. The signatures are pre-computed for all images in the image database.

III. IMAGE RETRIEVAL

The image retrieval based on query image uses signature comparisons. Given a query image, its signature is obtained using the steps outlined above. The Gaussian mixture associated with the image database is used for generating the signature of the query image. A similarity measure is used for comparing the signatures of the query image and the images in the database. An ordered list of images in the database is the result of the query.

There are several similarity measures that can be used for comparing two signatures. Kelly et al.[5] use the L_2 norm or the Euclidean distance given by (13) between the probability density functions.

$$D_E(I_1, I_2) = \left[\int_{\mathfrak{R}} (P_{I_1}(\mathbf{x}) - P_{I_2})^2 d\mathbf{x} \right]^{\frac{1}{2}}. \quad (13)$$

Another measure which can be used to compare two PDFs is the normalized inner product measure. This is given by

$$D_{nip}(I_1, I_2) = \frac{\int_{\mathfrak{R}} P_{I_1}(\mathbf{x})P_{I_2}(\mathbf{x})d\mathbf{x}}{\left[\int_{\mathfrak{R}} P_{I_1}^2(\mathbf{x})d\mathbf{x} \int_{\mathfrak{R}} P_{I_2}^2(\mathbf{x})d\mathbf{x} \right]^{\frac{1}{2}}}. \quad (14)$$

This associates a value between 0 and 1 as a measure of similarity between images. The normalized inner product measure is computed using

$$D_{nip}(I_1, I_2) = \left(\sum_{j=1}^K \sum_{l=1}^K w_j v_l \int_{\mathfrak{R}} G_j(\mathbf{x})G_l(\mathbf{x})d\mathbf{x} \right) \cdot \left(\sum_{j=1}^K w_j^2 \int_{\mathfrak{R}} G_j^2(\mathbf{x})d\mathbf{x} + 2 \sum_{j=1}^K \sum_{l=j+1}^K w_j w_l \int_{\mathfrak{R}} G_j(\mathbf{x})G_l(\mathbf{x})d\mathbf{x} \right)^{-\frac{1}{2}} \cdot \left(\sum_{j=1}^K v_j^2 \int_{\mathfrak{R}} G_j^2(\mathbf{x})d\mathbf{x} + 2 \sum_{j=1}^K \sum_{l=j+1}^K v_j v_l \int_{\mathfrak{R}} G_j(\mathbf{x})G_l(\mathbf{x})d\mathbf{x} \right)^{-\frac{1}{2}} \quad (15)$$

The integrals in (15) are computed using

$$\int_{\mathfrak{R}} G_j(\mathbf{x})G_l(\mathbf{x})d\mathbf{x} = (2\pi)^{\frac{N}{2}} |R_j + R_l|^{\frac{1}{2}} \exp \left[-\frac{1}{2}(c_1 + c_2) \right] \quad (16)$$

where c_1 and c_2 are given by:

$$\begin{aligned} c_1 &= \mathbf{m}_j^T R_j^{-1} \mathbf{m}_j + \mathbf{m}_l^T R_l^{-1} \mathbf{m}_l \\ c_2 &= -(\mathbf{m}_j^T R_j^{-1} + \mathbf{m}_l^T R_l^{-1})(R_j^{-1} + R_l^{-1})^{-1} \\ &\quad (\mathbf{m}_j^T R_j^{-1} + \mathbf{m}_l^T R_l^{-1})^T \end{aligned}$$

and

$$\int_{\mathfrak{R}} G_j^2(\mathbf{x})d\mathbf{x} = 2^{-N} \pi^{\frac{N}{2}} |R_j|^{\frac{1}{2}}. \quad (17)$$

In this paper, the Kullback-Leibler distance measure is used for comparing image signatures. The Kullback-Leibler distance measure between the images is given by

$$D_{KL}(I_1, I_2) = \int_{\mathfrak{R}} P_{I_1}(\mathbf{x}) \ln \frac{P_{I_1}(\mathbf{x})}{P_{I_2}(\mathbf{x})} d\mathbf{x}. \quad (18)$$

The results of image retrieval using this measure are compared with that of the normalized inner product.

The image signatures are probability density functions. The Kullback-Leibler measure gives the probabilistic distance of the image signature of image I_1 from that of image I_2 and hence, used as a measure for comparing images. The Kullback-Leibler measure is not a symmetric measure [4]. A symmetric measure is preferred for image retrieval. This is due to the requirement that, given two images I_1 and I_2 in the database, if I_2 is retrieved as a close match for a query using I_1 as an example, then I_1 should be retrieved as a close match for a query using I_2 as the sample image. A symmetric measure will ensure this. Hence, the symmetric divergence measure, $J(I_1, I_2)$, based on Kullback-Leibler measure is used instead.

$$J(I_1, I_2) = \int_{\mathfrak{R}} P_{I_1}(\mathbf{x}) \ln \frac{P_{I_1}(\mathbf{x})}{P_{I_2}(\mathbf{x})} d\mathbf{x} + \int_{\mathfrak{R}} P_{I_2}(\mathbf{x}) \ln \frac{P_{I_2}(\mathbf{x})}{P_{I_1}(\mathbf{x})} d\mathbf{x}. \quad (19)$$

Consider the computation of the symmetric divergence measure, (19), with $P_{I_1}(\mathbf{x})$ and $P_{I_2}(\mathbf{x})$ from (12). Since the image signature is a weighted sum of Gaussian distributions, and the Gaussian distributions are closed under convolutions [2, pg. 46], we approximate it to a Gaussian distribution. Hence, $P_{I_1}(\mathbf{x})$ and $P_{I_2}(\mathbf{x})$ are approximated to the Gaussian distributions $f_1(\mathbf{x}) = \mathcal{N}(\mu_1, \Sigma_1)$ and $f_2(\mathbf{x}) = \mathcal{N}(\mu_2, \Sigma_2)$, respectively, where

$$\begin{aligned} \mu_1 &= \sum_{j=1}^K w_j \mathbf{m}_j, & \Sigma_1 &= \sum_{j=1}^K w_j^2 R_j \\ \mu_2 &= \sum_{j=1}^K v_j \mathbf{m}_j, & \Sigma_2 &= \sum_{j=1}^K v_j^2 R_j. \end{aligned}$$

The probabilistic distance between $P_{I_1}(\mathbf{x})$ and $P_{I_2}(\mathbf{x})$ is approximated to the distance between $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$. A closed form representation [4] of this distance is

$$J(I_1 : I_2) = \frac{1}{2} \text{tr}[(\Sigma_1 - \Sigma_2)(\Sigma_2^{-1} - \Sigma_1^{-1})] + \frac{1}{2} \text{tr}[(\Sigma_1^{-1} + \Sigma_2^{-1})(\mu_1 - \mu_2)(\mu_1 - \mu_2)^T]. \quad (20)$$

A. Querying the Image Database

In query-based image retrieval from an image database, a user identifies a query image and is interested in retrieving “similar” images in the image database. The signature of the query image is computed as described in section II-D. The symmetric divergence measure given by (20) and the normalized inner product given by (15) are used for comparing the signature of the query image and the images in the database.

The output of the query is a ranked list of images in the database and their distance from the query image. The listing in the output is ordered in *decreasing order of similarity* with the most similar image listed first and the least similar image at the end.

IV. RESULTS AND OBSERVATIONS

The KODAK image database has 318 images. The indexing is based on HVC features and the data is clustered into 20 clusters. The normalized inner product measure and the symmetric divergence measure are used as a measure of similarity in image retrieval. Figure 1 shows the top 30 images retrieved from the database using the normalized inner product measure based on a query image. The query image is selected from the database itself and is retrieved as the first choice (the top left corner) in the result. Figure 2 shows the top 30 images retrieved from the KODAK image database using the symmetric Kullback-Leibler measure for the same query image.

In this particular query, these two measures have returned comparable results. However, the retrieval result using normalized inner product measure listed all images of the baby in the query image within top 6 images and the Kullback-Leibler distance measure returned these images within top 8 images. It should be noted here that the indexing and retrieval is based on the feature vectors in the HVC color space. No other semantic information is coded in representing these images. Hence images with similar color distribution are retrieved.

While more tests are required to discriminate between these two similarity measures, it is noted that the results of image retrieval using the Kullback-Leibler measure ranks other baby pictures in the database, specifically, at positions 15, 27, 28, 30 (within top 30). The normed inner product measure ranks those images at positions 108, 102, 95, 47 (respectively). A visual comparison of some of the results of image retrieval is not a sufficient test for comparing similarity measure for the purpose of image retrieval. A more robust method is required for this purpose.

V. CONCLUSIONS

This paper discussed image indexing based on pdf signatures and image retrieval using Kullback-Leibler measure and the normalized inner product measure as similarity measures. A probability density function is assigned as a signature of an image in the image database. Due to the complexity of the problem of pdf estimation, the signature was approximated to a linear combination of Gaussian distributions. The parameters of the Gaussian distributions, mean vector and covariance matrices, were estimated based on some clustering of the feature vectors of the images in the HVC color space. This method of generating image signatures was employed in indexing the KODAK image database.

The K -means clustering algorithm was used in computing the cluster centers in the feature space. This algorithm requires all the data to be available while generating the clusters. However, image databases are dynamic entities and hence a learning algorithm seems appropriate for clustering image features. Unsupervised Bayesian learning technique or other clustering techniques can be explored for this problem.

In this paper, the signature of an image is approximated to a Gaussian distribution and the distance between the two images is approximated to the symmetric divergence between their respective Gaussian signatures. Another method of approximating the image signatures is to invoke the minimum cross-entropy principle (MinxEnt). An image signature $P_I(\mathbf{x})$ can be approximated to a Gaussian distribution $f^{kl}(\mathbf{x})$ which minimizes the Kullback-Leibler distance

$$D_{kl}(P_I(\mathbf{x}) : f(\mathbf{x})) = \int_{\mathbb{R}} P_I(\mathbf{x}) \ln \frac{P_I(\mathbf{x})}{f(\mathbf{x})} d\mathbf{x}. \quad (21)$$

The distance between two image signatures, $P_{I_1}(\mathbf{x})$ and $P_{I_2}(\mathbf{x})$, is then given by the symmetric divergence between $f_1^{kl}(\mathbf{x})$ and $f_2^{kl}(\mathbf{x})$.

REFERENCES

- [1] R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern Classification (2nd Ed.)*. John Wiley & Sons, Inc., 1998, to be published.
- [2] W. Feller. *An Introduction to Probability Theory and its Applications*. John Wiley & Sons, Inc., New York, 1971.
- [3] Y. Gong and M. Sakauchi. Detection of regions matching specified chromatic features. *Computer Vision and Image Understanding*, pages 263–269, 1995.
- [4] J.N. Kapur and H.K. Kesavan. *Entropy Optimization Principles with Applications*. Academic Press, Inc., New York, 1992.
- [5] P. M. Kelly and M Cannon. Experience with candid: Comparison algorithm for navigating digital image databases. *In SPIE Vol 2368 Proceedings of the 23rd AIPR Workshop on Image and Information Systems: Applications and Opportunities*, pages 64–75, 1994.
- [6] P. M. Kelly, M Cannon, and J.E. Barros. Query by image example: the candid approach. *In SPIE Vol 2420 Storage and Retrieval for Image and Video Databases III*, pages 238–248, 1995.
- [7] J.N. Siddall. *Probabilistic Engineering Design: Principles and Applications*. Marcel Dekker, 1983.

