

## Robust Color Image Retrieval for the WWW

Bogdan Smolka \*

Polish-Japanese Institute of Information Technology  
Koszykowa 86, 02-008, Warsaw  
[bsmolka@pjwstk.edu.pl](mailto:bsmolka@pjwstk.edu.pl)

**Abstract.** The rapid growth of image archives increases the need for efficient and fast tools that can retrieve and search through large amount of visual data. In this paper an efficient method of extracting the image color content is proposed, which can serve as an image digital signature, allowing to efficiently index and retrieve large multimedia Internet based databases. The proposed method was applied using the images from the *WEBMUSEUM* Internet database containing the collection of images of fine arts. The results show that the new method of image color representation is robust to image resizing and compression and can be incorporated into existing web-based image retrieval systems.

### 1 Introduction

Successful queries on large, distributed databases cannot rely on textual information only and therefore color image indexing is one of the most important methods used for automatic content based retrieval. In this paper we focus on the image indexing, based on the global color distribution, which is applied for cases when the user provides a sample image for the query.

The majority of the systems exploiting the image color information work using various kinds of color histograms. However the histogram based approach has many drawbacks, as the histogram representation is sensitive to illumination changes, image resizing through interpolation and compression induced artifacts. Therefore, in this paper we propose a nonparametric approach to the problem of the estimation of the distribution of image colors.

### 2 Color Histograms

Color indexing is a process through which the images in a database are retrieved on the basis of their color content. The indexing process must enable the automatic extraction of features, efficient assigning of digital signatures to images and effective retrieval of images within a database.

In order for an image retrieval system to retrieve images that are visually similar to the given query, a proper representation of the visual features is needed and a measure that can determine the similarity between a given query and the images from a database

---

\* This research has been supported by a grant No PJ/B/01/2004 from the Polish-Japanese Institute of Information Technology

set has to be chosen. Assuming that no textual information about the image content are given, image features such as color [1, 2, 3], texture [4] and shape [5, 6] are commonly used.

These features are dependent on illumination, shading, resizing manipulations and compression induced artifacts. Thus, the visual appearance of an image is better described by the distribution of features, rather than by individual feature vectors.

Color feature has proven to be efficient in discriminating between relevant and non-relevant images. One of the widely used tools for image retrieval is the color histogram, which describes the distribution of colors in an image using a specific color space. The colors of an image are mapped into a discrete color space containing  $m$  colors. In this way, a color histogram is an  $m$ -dimensional vector, whose elements represent the number of pixels of a given color in an image.

In this paper we use the RGB color space, which although not perceptually uniform, is the most commonly used, primarily to retain compatibility with computer display systems. Let us assume a color image  $\mathbf{F}$  of size  $n_1 \times n_2 = N$ , composed of three RGB channels  $\mathbf{F} = \{F_{i,j}^R, F_{i,j}^G, F_{i,j}^B\}$ ,  $i = 1, \dots, n_1, j = 1, \dots, n_2$ .

An image histogram  $H$  in the RGB color space is the simplest approximation of the density function of the image RGB channels intensities

$$H(\rho, \gamma, \beta) = \# \{F_{i,j}^R = \rho, F_{i,j}^G = \gamma, F_{i,j}^B = \beta\} / N, \quad (1)$$

where  $N$  is the total number of image pixels, and  $\#$  denotes the number of pixels with color channel values  $\{\rho, \gamma, \beta\}$ .

For the analysis of colors, which is independent of image brightness, it is convenient to transform the RGB values into normalized components  $r, g, b$  defined as:  $r = R/I$ ,  $g = G/I$ ,  $b = B/I$ ,  $I = R + G + B$ , where  $R, G, B \in [0, 255]$ . The normalized color values can be expressed using only  $r$  and  $g$  values as  $b = 1 - r - g$  and the normalization makes the  $r, g$  variables non-dependent on the brightness  $I$ .

Using the normalized  $rg$  reduced color space, we can map the color pixel on a two-dimensional plane and obtain a two-dimensional discrete histogram

$$\begin{aligned} \Phi(x, y) &= \# \{ \text{int}(MF_{i,j}^R/I_{i,j}) = x, \text{int}(MF_{i,j}^G/I_{i,j}) = y \} / N \\ &= \# \{ \text{int}(Mr_{i,j}) = x, \text{int}(Mg_{i,j}) = y \} / N, \end{aligned} \quad (2)$$

$x, y = 0, \dots, M$ , where  $M + 1$  is the dimension of the 2-dimensional histogram, (for true-color images  $M = 255$ ). The likeness between two images is often expressed through the similarity of their color histograms. One of the most popular ways to measure the similarity between two histograms is the Minkowski distance or the histogram intersection, which were also used in this work, [1, 7].

### 3 Nonparametric Color Distribution

The drawback of the histogram representation is that the shape of the histogram strongly depends on the method used for lossy image representation and on the number of image pixels, as for small image sizes there is too few points to build the 3-dimensional color histogram, which makes that the comparison of histograms is inapplicable.

To alleviate the problems, we propose in this paper to estimate the color distribution not through the discrete histogram, but to use a smooth nonparametric estimate, based on the concept of nonparametric density estimation, [8, 9]. In this formulation, the similarity measure between two estimates of the color distribution will be expressed as the distance between two surfaces of the two-dimensional kernel density estimation in the normalized  $rg$  color space.

*Density Estimation* describes the process of modelling the probability density function  $f(x)$  of a given sequence of sample values drawn from an unknown density distribution. The simplest form of density estimation is the histogram, however its main disadvantage is its strong dependence on the chosen bin-width and the origin of the histogram grid.

Kernel Density Estimation, avoids this disadvantage by placing a kernel function on every sample value in the sample space and then summing the values of all functions for every point in the sample space, (Figs. 1, 2, 3). This results in a smooth density estimates that are not affected by an arbitrarily chosen partition of the sample space. The multivariate kernel density estimator in the  $q$ -dimensional case is defined as

$$\hat{f}_{\mathbf{h}}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_1 \cdots h_q} \mathcal{K} \left( \frac{x_{i1} - x_1}{h_1}, \dots, \frac{x_{iq} - x_q}{h_q} \right), \quad (3)$$

with  $\mathcal{K}$  denoting a multidimensional kernel function  $\mathcal{K}: \mathbb{R}^q \rightarrow \mathbb{R}$  and  $h_1, \dots, h_q$  denoting bandwidths for each dimension and  $n$  is the number of samples in the sliding window. A common approach to build multidimensional kernel functions is to use a *product kernel*  $\mathcal{K}(u_1, \dots, u_q) = \prod_{i=1}^q K(u_i)$ , where  $K$  is a one-dimensional kernel function. Intuitively, the kernel function determines the shape of the 'bumps' placed around the sample values and the bandwidths  $h_1, \dots, h_q$  their width in each dimension. If bandwidth is the same in all dimensions, multivariate radial-symmetric kernel functions can be used,

$$\hat{f}_h(\mathbf{x}) = \frac{1}{nh^q} \sum_{i=1}^n K \left( \frac{\|\mathbf{x}_i - \mathbf{x}\|}{h} \right). \quad (4)$$

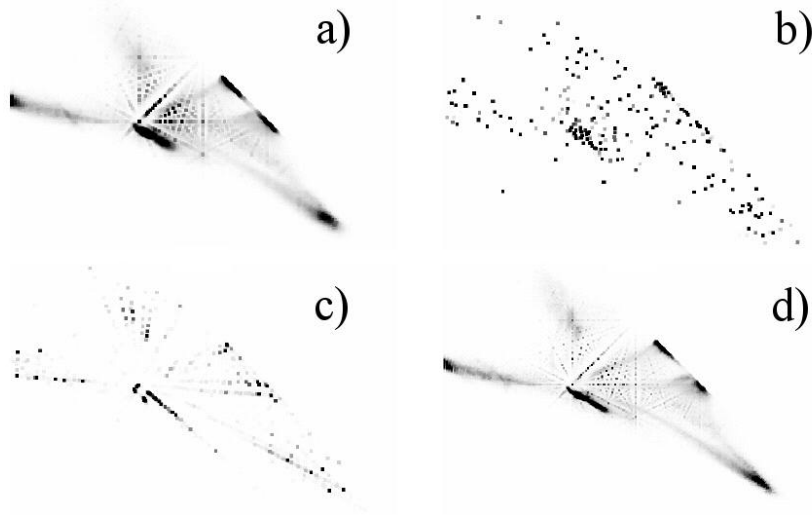
The shape of the approximated density function depends heavily on the bandwidth chosen for the density estimation. Small values of  $h$  lead to spiky density estimates showing spurious features. On the other hand too big values of  $h$  produce over-smoothed estimates that hide structural features.

The unknown density function is commonly assumed to be the normal distribution and choosing the Gaussian kernel for  $K$ , the optimal bandwidth in the one-dimensional case is:  $h_{opt} = 1.06\hat{\sigma}n^{-\frac{1}{5}}$ , where  $\hat{\sigma}$  denotes the standard deviation, and for the  $q$ -dimensional case, [8]

$$h_{opt} = (4/(q+2))^{\frac{1}{q+4}} \hat{\sigma} n^{-\frac{1}{q+4}}. \quad (5)$$

In this paper we use the  $rg$  color space, so  $q = 2$  and we used the Gaussian kernel, although we have obtained similar results using other kernel shapes commonly used in nonparametric density estimation.

Using the kernel based estimation, a smooth estimate of the color distribution is obtained as shown in Figs. 2, 3, 4. As can be seen in Fig. 4, the density distribution is



**Fig. 1.** Influence of the compression methods on the color distribution in the normalized  $rg$  color space: a) test image PILLS (Fig. 2) of size 512x512, b) PILLS in GIF format, c) PILLS in JPEG format (compression ratio 78), d) PILLS in JPEG2000 (compression ratio 120).

insensitive to resizing and lossy image coding, which are the basic operations performed when preparing large Internet multimedia databases. This distribution can be used for the image retrieval purposes, as it can serve as an image signature, as depicted in Fig. 3, which shows the  $rg$  nonparametric distributions of some well known color test images.

## 4 Results

To evaluate the efficiency of the proposed color density estimation, we used as the testbed the collection of 3000 images comprising the well-known *WEBMUSEUM* Internet database, ([www.ibiblio.org/wm/](http://www.ibiblio.org/wm/)) and the collection of 10 000 low resolution web-crawled jpeg images from the database of J.Z. Wang, ([wang.ist.psu.edu/docs/related/](http://wang.ist.psu.edu/docs/related/)).

The first database contains a collection of about 3000 images of fine arts of various famous artists. Each image is coded in JPEG of moderate compression ratio, (the blocking artifacts are hardly visible) with width or height of about 1000 pixels. Each image is accompanied by a thumbnail of width or size of 100 pixels, also compressed with JPEG.

From the database, the image of the painting "Starry Night" of V. van Gogh was chosen as the query image, (see Fig. 5 left). Using the kernel density estimation, we used the Euclidean distance between the surfaces of the color distribution estimate as the similarity measure and ordered the retrieved images according to the distance values. The results are very promising, as the second image, most similar to the query, was another painting of van Gogh, "Road with Cypress and Star", (Fig. 5, second image in the row).

In the second experiment, (Fig. 6) we used the thumbnail of the *Starry Night* image as a query image. Although this picture is small ( $122 \times 100$ ) and heavily jpeged, the proposed scheme was able to find the image of full resolution, (first image in the ordered sequence) and the majority of images retrieved using the full resolution images. Very similar results were obtained using the histogram intersection method, so as expected the two methods of similarity evaluation yield comparable results.

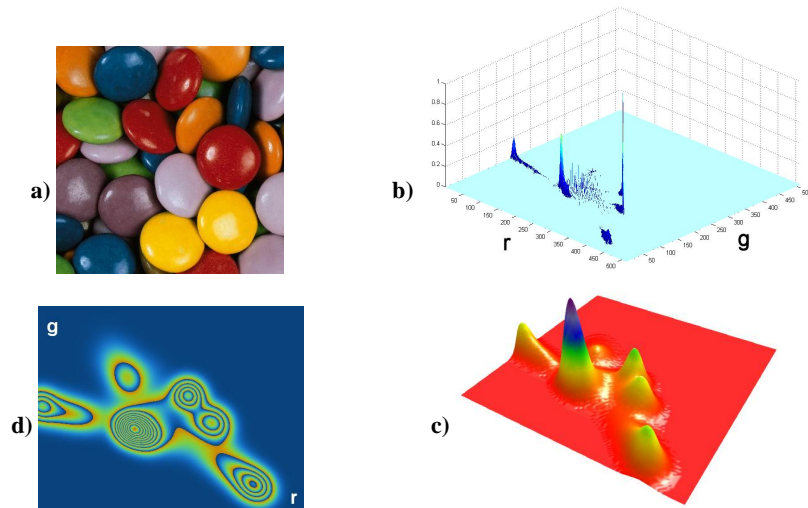
Fig. 7 shows the most similar images from the WEBMUSEUM database to the well known test image LENA. Surprisingly the most similar image was the "Girl with a Pearl Earring" of Vermeer. Additionally, Fig. 8 shows the most similar images to the "Starry Night" from the database of Wang.

## 5 Conclusions

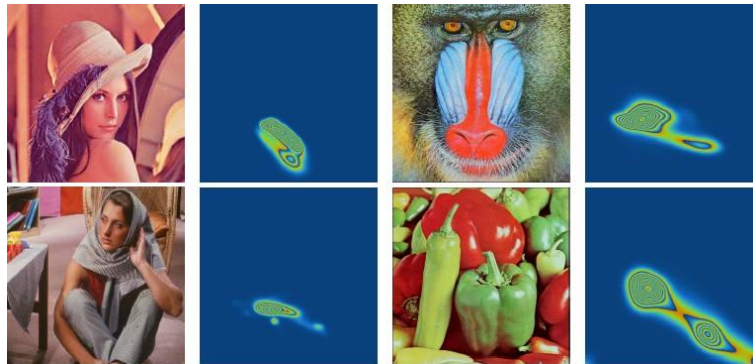
In this paper we proposed a robust way of color density estimation. To enable fast retrieval of large databases we used the normalized  $rg$  color space. The experiments show that the method of nonparametric density estimation is insensitive to image compression and resizing. This makes the proposed framework interesting for image retrieval applications. Especially, the ability to retrieve images using a heavily distorted thumbnail is interesting, as it enables extremely fast retrieval of large databases.

## References

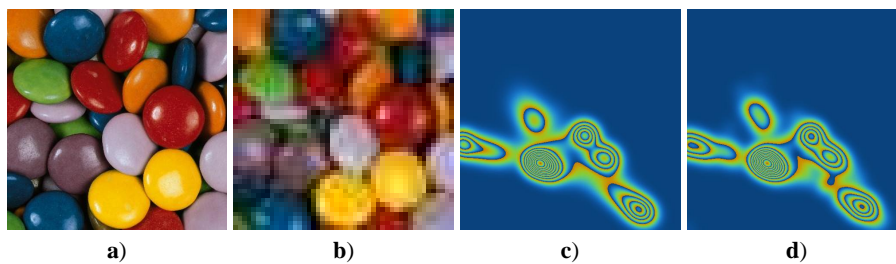
1. Swain M., D. Ballard, Color indexing, International Journal of Computer Vision, 7, 1, 11-32, 1991.
2. M. Stricker, M. Orengo, Similarity of color images, in SPIE Conference on Storage and Retrieval for Image and Video Databases III, 2420, 381-392, February 1995.
3. X. Wan, C.C.J. Kuo, Color distribution analysis and quantization for image retrieval, Proceedings of SPIE, Vol. 2670, February 1996.
4. B.S. Manjunath, W.Y. Ma, Texture features for browsing and retrieval of image data, IEEE Transactions on Pattern Analysis and Machine Intelligence, 18, 8, 837-842, 1996.
5. J.E. Gary, R. Mehrotra, Similar shape retrieval using a structural feature index, Information Systems, 18, 7, 525-537, October 1990.
6. A.K. Jain and A. Vailaya. Image retrieval using color and shape. Pattern Recognition, 29, 8, 1233-1244, 1996.
7. B.S. Manjunath, J.R. Ohm, V.V. Vasudevan, A. Yamada, Color and texture descriptors, IEEE Trans. CSVT, 11, 6, 703-715, June 2001.
8. D.W. Scott, "Multivariate Density Estimation", New York, John Wiley, 1992.
9. B.W. Silverman, "Density Estimation for Statistics and Data Analysis", London, Chapman and Hall, 1986.



**Fig. 2.** Illustration of the nonparametric probability density estimation: **a)** test image PILLS and its histogram in the rg color space (**b**), **c)** and **d)** present the visualization of the smooth kernel based estimation.



**Fig. 3.** Examples of density estimation of the color distribution in the rg space.



**Fig. 4.** Robustness of density estimation to image resizing and compression: **a)** test image of size 512x512, **b)** resized with bilinear interpolation (48x48) and then compressed with JPEG, (compression ratio 8.3), **c)**, **d)** pseudo-color representation of the normalized color densities of **a)** and **b)** respectively).





**Fig. 5.** Results for the query for images similar to the van Gogh "Starry Night" painting (first image) from the WEBMUSEUM database.



**Fig. 6.** Results for the query for images similar to the thumbnail of size (122 x 100) of the van Gogh "Starry Night" painting, from the WEBMUSEUM.



**Fig. 7.** Results for the query for images from the WEBMUSEUM database similar to the LENA color test image.



**Fig. 8.** Results for the query for images from the database of Wang similar to the full resolution "Starry Night" image.