

Biomarker Detection in Whole Slide Imaging based on Statistical Color Models

Jie Shu¹, Guoping Qiu^{1*}, Mohammad Ilyas^{2,3} and Philip Kaye³

¹ School of Computer Science, University of Nottingham, UK

² School of Molecular Medical Sciences, University of Nottingham, UK

³ Queens Medical Center NHS Trust, Nottingham, UK

Abstract. This paper presents a technique for immunostaining biomarker detection in digital slides. We treat immunostaining detection as a color image analysis problem and build statistical color models using a large number of labeled positive and negative immunostaining pixels. We have implemented the statistical models in different color spaces and show that the opponent chromaticity signals effectively characterize the color distributions of the immunostaining biomarkers and that the luminance is an unreliable and distractive signal. We have applied the technique to the detection of positive P53 immunostaining in digital slides of oesophagitis and colorectal biopsies. We present experimental results and show that the technique can achieve a biomarker detection rate of over 98% with 5% false positives.

Keywords: Immunostaining, whole slide imaging, digital pathology, color image analysis, statistical models, ImageJ.

1 Introduction

Digital slide technology is changing the ways in which pathologists have been practicing disease diagnosis and prognosis for over 100 years. Instead of totally relying on inspecting glass slides through the microscopes, pathologists can now exploit computer and computerized image analysis technology to aid diagnosis. Converting glass slides into digital slides (also called whole slide imaging) has many potential benefits which include, automating some of the diagnosis procedures thus enhancing efficiency; and reducing inter-observer discrepancies thus improving reproducibility and diagnosis accuracy.

In essence, digital pathology consists of scanning the glass slides into digital images, displaying the digital slides on devices such as LCD monitors, managing and archiving digital slide files on networked computer systems, and using image analysis software tools to analyze the digital slides for making diagnosis and prognosis. In this paper, we present a study of developing image analysis and classification techniques to assist pathologists interpret digital slides.

The promising potential of whole slide imaging has stimulated much recent research and commercial activities in digital pathology [1]. Despite much interest and

* Send all correspondence to: qiu@cs.nott.ac.uk

huge potential, systematical application of computer vision and image analysis algorithms to histopathology has only just begun. Although tasks such as identifying positive immunostaining pixels (cells) is similar to that of object detection in computer vision, applications of existing image analysis will still be very challenging. One of the reasons is that many computer vision algorithms are still at a research stage and not yet sufficiently stable to work robustly on challenging tasks such as accurately identifying cells and cell structures. One of the biggest challenges in computer vision is the huge variability of the input objects and environments. For example, in object recognition, the image of the object changes with lighting conditions and the camera's positions; the complexity of the backgrounds makes it hard to separate the objects from the background clutter, etc. This is an ongoing research field and algorithms that can robustly tackle the variability of the inputs are current research topics in computer vision. Many medical image analysis problems face the same and even greater challenges as in computer vision.

For the case of digital slide analysis, the challenges include robustly identifying positive immunostaining, cell segmentation and tissue structure recognition. From a pure image analysis perspective, the challenges include that, the color distributions will vary with tissue samples, diseases, the concentration of chemicals, the experience of the lab technicians and the scanning instruments; the shape and sizes of the cell will depend on the direction of section, the thickness of the sample, and the tissue morphological structures will vary greatly from sample to sample. These factors on top of the immaturity of image analysis technology have added challenges to tissue image analysis.

One of the first steps in digital slide analysis is the identification of positive immunostaining in the digital slides. In computer vision terms, this is a color image analysis problem. In the case of P53 immunostaining, potentially positive cancer cells will appear brown. The objective of detecting positive stains is therefore to identify brown pixels. However, this seemingly simple visual recognition task is in fact not so simple in computer analysis. When an experienced pathologist sees positive staining, he/she may find it relatively easy to identify them. However, there are a variety of different shades of brown colors within one single slide, and to complicate the matter, different slides will have a different range of shades of brown. The question is: what browns are positive browns? Although an experienced pathologist can identify them easily when he/she sees them, but it is very hard for him/her to describe it clearly and concisely, and it is even harder to convert the pathologists descriptions into numerical formulas – the forms that a computer analysis algorithm rely on. To add to the complication, different pathologists may regard different shades of brown as positive/negative thus causing inter observer discrepancies. This is rather like a face recognition problem. Humans can recognize faces very easily, and yet we still do not know how our brain does this. This is a classic “I will know it when I see it” scenario.

In this paper, we tackle automatic immunostaining detection in digital slide using statistical models. We borrow techniques developed in color image analysis and use very large collection of digital slide pixels to build statistical color models for positive and negative immunostaining pixels. Based on the statistical models of the positive and negative staining pixels, we use a maximum likelihood classifier to automatically classify pixels in digital slides into positive and negative staining pixels. We present

experimental results on digital slides of oesophagitis and colorectal biopsies and demonstrate the effectiveness of our method.

2 The Problem

Immunostaining for the evaluation of antigen expression is a standard procedure for the diagnosis and prognosis of cancer and other diseases [2]. Increasingly, monitoring changes in targeted antigens has been used to measure cancer treatment response and disease progression. Routine diagnosis is predominantly performed by visual inspection and scoring. Many factors affect the diagnosis accuracy and reliability of immunostaining and amongst them are the problems associated with manual inspecting and diagnosis which are heavily relied on the experience of the observers and there are often large discrepancies between observers. Computerized image analysis offers the potential of automating the diagnosis process and reducing inter observer variability.

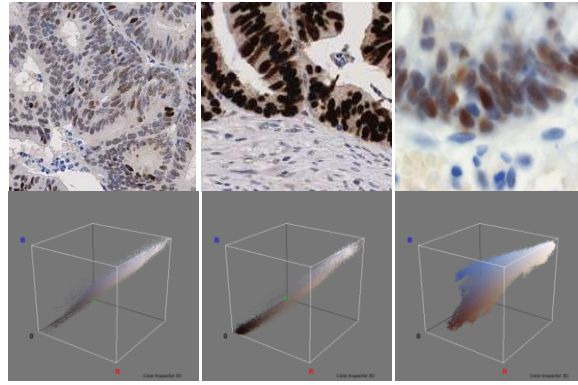


Fig. 1. The color distributions of different tissue samples stained with the same immunostaining (P53) vary significantly. Top row: image samples. Bottom row: plots of the color distribution in the RGB color space.

The problem that we address in this paper is immunostaining biomarker detection in digital slides using color image analysis. Specifically, we want to develop image analysis algorithms for the automatic detection of positive P53 stains in whole slide imaging of biopsies. Visually, positive P53 stains appear as brown colors. However, depending on sample preparation, the types of tissues and the types of disease and data acquisition procedures and equipments, there are very different shades of brown which make automatic detection of positive staining difficult. Fig. 1 show three samples of biopsies stained with P53 immunostaining. Visually, it is seen that there are a range of different shades of brown pixels within a slide and they vary hugely across different samples.

Another problem associated with whole slide imaging diagnosis is the huge image size. A single slide can be as large as 50,000 x 50,000 pixels. Sieving through such an enormous amount of pixels manually is a huge task. What is desired is an automatic

solution which will guide pathologists to inspecting region of potential interest rather than having to go through every pixel of the image.

3 Statistical Color Models for Immunostaining Detection

Color has been extensively used in computer vision for object detection. There is a large body of color image analysis literature. Color has also been used for immunostaining detection [3, 4]. In [8] color is converted to grayscale for cell segmentation. In [6], a simple condition $\text{Red} > \text{Blue}$ has been used to identify staining, [7] used 11 different color spaces and machine learning for cancer classification, and [9] quantitatively studied immunostaining using the CMYK color model.

In this paper, we take a statistical approach. We collect a very large number of labeled pixels and build the statistical models for the positive and negative staining pixels, and use a maximum likelihood classifier to classify the immunostaining pixels. Based on the observation that it is the color information that will play the most significant role in the identification of the immunostaining [7], we have implemented solutions using only the chromaticity signals. We shown that using the chromaticity signals not only reduces the model complexity but also it is a more effective way of characterizing and modeling the immunostaining colors.

3.1 Semi-automatic sample labeling and collection

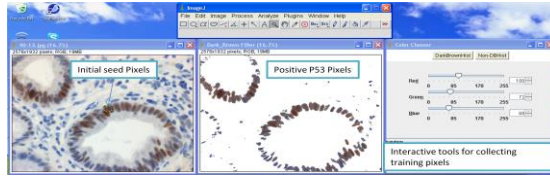


Fig. 2. Screen shot of the semi-automatic tool for collecting training examples. User chooses a set of initial seed pixels of interest; the positive pixels are then visualized; an interactive tool can then be used to refine the selection.

As in building any automatic machine classification algorithms, the first thing we need is training samples. We have built a semi-automatic tool using ImageJ to collect positive staining pixels and background pixels. The tool is illustrated in Fig.2. This is a simple tool, users can use the sliding bars to set the combination of color axis and visually determine the positive staining pixels and background. These manually labeled pixels are then used to build the color statistical models.

3.2 Statistical Color Model Construction

As described in [12], statistical color models can be constructed using the histograms of the color images. The probability of a color, represented in a color space rgb can be

obtained by first quantizing the color into color bins [rgb] and then by counting the number of pixels that falling into each of the bins, i.e.

$$Prob(rgb) = \frac{\#[rgb]}{\# \text{ Total Pixel}} \quad (1)$$

To construct a statistical model for the positive and negative staining pixels, we will use the labeled positive and negative pixels respectively, i.e.

$$Prob(rgb|S) = \frac{\#S[rgb]}{N_s} \quad Prob(rgb|\bar{S}) = \frac{\#\bar{S}[rgb]}{N_{\bar{s}}} \quad (2)$$

where S represents positive staining and \bar{S} represents negative staining respectively and N_s is the total number of positive and $N_{\bar{s}}$ is the total number of negative pixels.

It is known in the color image analysis literature, in general, the Red, Green and Blue (rgb) color space is not suited for image analysis. One of the reasons is that chromaticity information and brightness (luminance) information are mixed together in this color space and it is often desirable to be able to process the chromatic signals and the luminance signal separately. It is therefore a common practice to separate the chromaticity signal from the luminance signal. In immunostaining detection, it is the chromaticity signal or the color spectral that is of interest rather than the absolute brightness. The chromaticity signals encode the spectral information of the immunochemical and therefore can be used to detect positive staining. From a computational perspective, using a 2D chromaticity space will make it easier to model the probability density function (curse of dimensionality problem).

There are many color models, in this paper, two opponent color models often used in computer vision literature have been tested. The first is the red-green and blue-yellow opponent space, and the second is the Cb and Cr space of YCbCr model.

The red-green (rg) and blue-yellow (by) chromaticity signals are derived from the original RGB input as follows

$$r = \frac{R}{R+G+B} \quad g = \frac{G}{R+G+B} \quad b = \frac{B}{R+G+B} \quad rg = r - g \quad by = \frac{r+g}{2} - b; \quad (3)$$

The Cb and Cr chromaticity signals are derived from the original RGB space as follows;

$$\begin{aligned} Cb &= -0.1687 * R - 0.3313 * G + 0.5 * B \\ Cr &= 0.5 * R - 0.4187 * G - 0.0813 * B \end{aligned} \quad (4)$$

The color statistical models in the chromaticity space can now be constructed as

$$Prob((rg, by)|S) = \frac{\#S[rg, by]}{N_s} \quad Prob((rg, by)|\bar{S}) = \frac{\#\bar{S}[rg, by]}{N_{\bar{s}}} \quad (5)$$

$$Prob((Cb, Cr)|S) = \frac{\#S[(Cb, Cr)]}{N_s} \quad Prob((Cb, Cr)|\bar{S}) = \frac{\#\bar{S}[(Cb, Cr)]}{N_{\bar{s}}} \quad (6)$$

3.3 Maximum Likelihood Biomarker Detection

Based on the statistical models of the positive stain and negative stain pixels, we can use the likelihood ratio approach to build the classifier [10]. A pixel is classified as positively stained if

$$\frac{\text{Prob}(\text{rgb}|\text{S})}{\text{Prob}(\text{rgb}|\bar{\text{S}})} \geq \theta \quad (7)$$

where $0 \leq \theta \leq 1$ is the threshold. The value of θ trades-off correction detection and false positive which is the most important property of the classifier. One possible way to determine the threshold value can be determined as follows [11]

$$\theta = \frac{c_p \text{Prob}(\text{S})}{c_n \text{Prob}(\bar{\text{S}})} \quad (8)$$

where c_p and c_n are the application dependent costs associated with false positive and false negative. One possible way to compute $\text{Prob}(\text{S})$ is

$$\text{Prob}(\text{S}) = \frac{N_S}{N_S + N_{\bar{\text{S}}}} \quad (9)$$

4 Experimental Results

We have tested the technique on two sets of digital slides of oesophagitis and colorectal biopsies. The slides were scanned using a Hamamatsu scanner. We first used a semi-automatic tool to manually label positively stained pixels as described in section 3.1. The data for building the statistical color models include 4,429,970 positive immunostaining pixels and 393,553,963 background pixels from the biopsy samples. The models are then tested on another set of images consisting of 4,400,900 positive and 393,583,033 background pixels.

Table 1. The number of bins occupied by the positively stained pixels (S), background pixels ($\bar{\text{S}}$) and overlapping (OL) bins in the color histograms of 4096, 16384, 65536 bins (rg/by and Cb/Cr) and 262144, 2097152 and 16777216 bins (RGB and YCbCr)

	64 bins in each axis			128 bins in each axis			256 bins in each axis		
	S	$\bar{\text{S}}$	OL	S	$\bar{\text{S}}$	OL	S	$\bar{\text{S}}$	OL
rg/by	449	588	202	1296	1368	403	3149	3099	496
Cb/Cr	56	155	18	181	516	49	514	1308	129
RGB	6553	13735	2082	35915	79150	9001	61372	139005	14354
YCbCr	1991	4122	698	11038	23182	2972	44033	97609	10549

As expected, the color distribution occupies a relatively small part of the color space. Table 1 shows the number of non-empty bins for different sizes of the color histograms in 4 different color spaces including 2 opponent chromaticity spaces the RGB and YCbCr spaces. Also shown are the number of overlapping bins, i.e., the bins that are occupied by both the positive and negative stained pixels. It is interesting to observe that CbCr chromaticity space has the least number of overlapping bins and experimental results confirm that this space gave the best performances.

Fig. 3 shows the ROC curves of models of different color spaces of histograms of different bins. It is seen that although all models achieved very good results, the opponent chromaticity models achieved better results than models including both chromaticity and luminance. This indicates that chromaticity is sufficient and the luminance is a distraction in terms building the model. As mentioned before, the 2D

chromaticity signals make the model simpler, fast to compute and demand less memory.

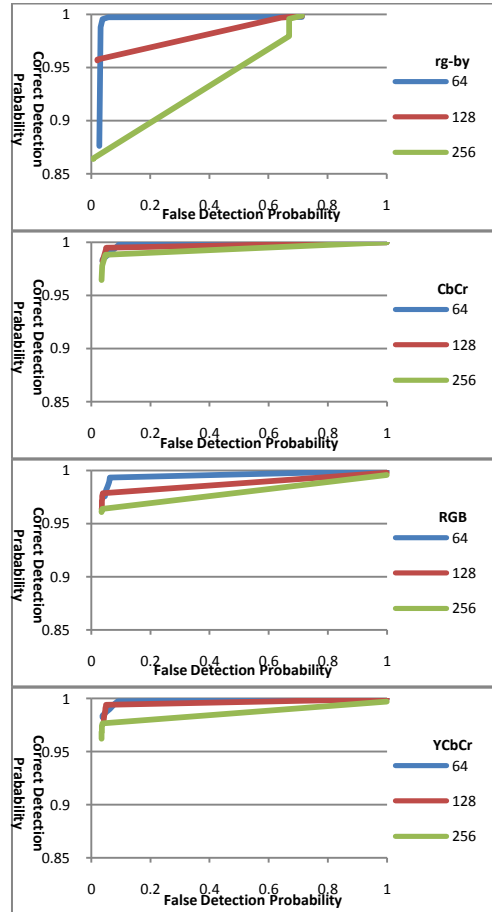


Fig. 3. ROC curves of immunostaining detection.

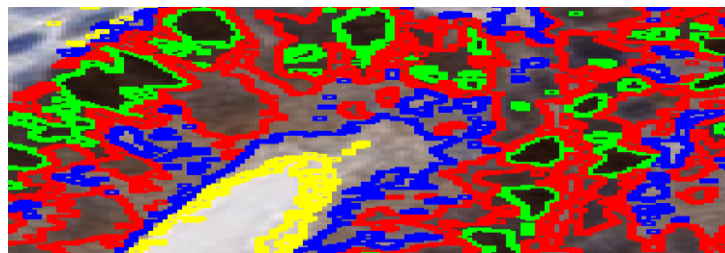


Fig. 4. An example of immunostaining detection using statistical color models in the CbCr chromaticity space. The region circled by different colors represent different shades of brown (from dark to light, green, red, blue, yellow) corresponding to Y (luminance) intensity of 0 – 63, 64 – 127, 128 – 191 and 192 – 255. This is useful for pathologists to inspect different shades of brown (positive staining colors)

Fig. 4 shows a detection example. Because we have separated chromaticity from luminance, we can simply use the luminance values to separate different shades of brown very easily, which will be convenient for pathologists.

5 Concluding Remarks

In this paper, we have developed a method that uses color statistical models for the detection of immunostaining in whole slide imaging. We have shown that using opponent chromaticity color signals is sufficient in building the statistical model for detecting the biomarkers in the digital slides. The advantages of using the chromaticity are that we can build a simpler model and separating chromaticity from the luminance also helps pathologists to investigate different shades of brown more easily. We believe the method will be a useful tool to facilitate pathologists to find regions of interest for diagnosis and prognosis.

Finally, our technique has been implemented as an ImageJ plugin and will be made available to the public in due course.

References

1. M. May, "A better lens on disease", *Scientific America*, May 2010, pp. 56 – 59
2. J. E. Anderson, et al, "Methods and biomarkers for the diagnosis and prognosis of cancer and other diseases: towards personalized medicine", *Drug Resist Updates*, 9(4-5):198-210, 2006
3. R. A. Walker RA, "Quantification of immunohistochemistry—issues concerning methods, utility and semiquantitative assessment I", *Histopathology*, 2006; 49:406 – 10
4. C. R. Taylor CR and R. M. Levenson, "Quantification of immunohistochemistry—issues concerning methods, utility and semiquantitative assessment II", *Histopathology*, 2006; 49:411–24.
5. Matos et al, "Immunohistochemistry as an important tool in biomarkers detection and clinical practice", *Biomarker Insight*, *Biomarker Insights* 2010:5 9–20, available from <http://www.la-press.com>
6. A S. Joshi et al, "Semi-Automated Imaging System to Quantitate Her-2/neu Membrane Receptor Immunoreactivity in Human Breast Cancer", *Cytometry Part A* 71A:273–285 (2007)
7. Mutlu Mete, Umit Topaloglu, "Statistical comparison of color model-classifier pairs in hematoxylin and eosin stained histological images", *Proceedings of the 6th Annual IEEE conference on Computational Intelligence in Bioinformatics and Computational Biology*, 2009, pp. 284-291
8. Mao KZ, Zhao P, Tan PH, Supervised learning-based cell image segmentation for p53 immunohistochemistry, *IEEE Trans Biomed Eng.* 2006 Jun;53(6):1153-63
9. Nhu-An Pham et al, "Quantitative image analysis of immunohistochemical stains using a CMYK color model", *Diagnostic Pathology*, published 27 Feb 2007, volume 2.
10. Keinosuke Fukunaga, "Introduction to Statistical Pattern Recognition", Academic Press, 1972
11. Harry L. Van Trees. "Detection, Estimation, and Modulation Theory", volume I. Wiley, 1968.
12. Jones, M. J. and Rehg, J. M. 2002. Statistical color models with application to skin detection. *Int. J. Computer Vision* 46, 1 (Jan. 2002), 81-96