

Advanced Machine Learning Technique for Classification of Histopathological Images

Miss. Pooja Mali¹

PG Student

Department of Computer Engineering JSPM Narhe Technical Campus, Narhe Pune-411041

Prof. Vilas. S. Gaikwad²

Assistant Professor Department of Computer Engineering JSPM Narhe Technical Campus, Narhe

Pune-411041

Abstract—The classification of breast cancer has been the subject of interest in the fields of healthcare and bioinformatics, because it is the second main reason of cancer-related deaths in women. Breast cancer can be analyzed using a biopsy where tissue is eliminated and studied under microscope. The identification of problem is based on the qualification and experienced of the histopathologists, who will attention for abnormal cells. However, if the histopathologist is not well-trained or experienced, this may lead to wrong diagnosis. With the recent proposition in image processing and machine learning domain, there is an interest in experiment to develop a strong pattern recognition based framework to improve the quality of diagnosis. In this work, we will use the image feature extraction approach and machine learning approach for the classification of breast cancer using histology images into benign and malignant. Using Histopathological image we can preprocess this image after that apply feature extraction and classify the final result using SVM and Naive Bayes Classification techniques.

Index Terms—Histopathological image classification, breast cancer diagnose, feature extraction, SVM classification, Naive Bayes Classification;

they cannot be used to find or identified whether the area is cancerous. The biopsy, where a tissue is gives as input and processed under a microscope to see if cancer is present, is the only sure way to find if an area is cancerous. After completing the biopsy, the identification of problem will be based on the qualification of the histopathologists, who will analyze the tissue under a microscope, looking for exceptional or cancerous cells. The histology images allow us to differentiate the cell nuclei types and their flowchart according to a specific pattern. Histopathologists particularly examine the consistency of cell shapes and tissue distributions and decided the cancerous regions and malignancy degree. If the histopathologists are not well-trained, this may lead to an incorrect identification of problem. Also, there is a lack of specialists, which maintain the tissue sample on hold for up to two months. There is also the issue of reproducibility, as histopathology is a subjective science. This is right especially between non-specialized pathologists, where we can get a different identification of problem on the same sample. Therefore, there is an insistent demand for computer-assisted identification of problem.

I. INTRODUCTION

A. BACKGROUND

Breast cancer is the most common and dangerous intrusive cancer in women and the second main effect of cancer death in women, after lung cancer. The International Agency for Research on Cancer (IARC), which is part of the World Health Organization (WHO), the numbers of deaths reasoned by cancer in the year of 2012 only come to around 8.2 million. The number of new cases is expected to growth to more than 27 million by 2030.

Finding breast cancer quick and getting state-of-the-art cancer treatment are the key plan of action to avoid deaths from breast cancer. In existing, it is a widely-used way to identification of breast cancer by identifying hematoxylin and eosin (HE) stained histological slide preparations that are checked under a high powered microscope of the changed area of the breast. In medical practice, classification of breast cancer biopsy result into different plans (e.g. cancerous and noncancerous) is manually driven by experienced pathologists. Come out machine learning approaches and enlarging image volume developed automatic system for breast cancer classification possible and can help pathologists to obtain precise identification of problem more efficient. Breast cancer can be find or identified using medical images testing using histology

and radiology images. The radiology images search can help to find the areas where the difference is located. However,

B. MOTIVATION

Breast cancer can be identified using a biopsy where tissue is removed and studied under microscope. The identification of problem is based on the qualification and experienced of the histopathologist, who will look for abnormal cells. However, if the histopathologist is not well-trained or experienced, this may lead to wrong identification of problem. The recent proposition in image processing and machine learning domain, there is an interest try to develop a reliable pattern recognition based approach to improve the quality of identification of problem.

C. OBJECTIVES

1. To classify the breast cancer histology images into benign and malignant
2. To work on histopathological image dataset for breast cancer classification
3. To developed features-based classification methods.

II. REVIEW OF LITERATURE

Breast cancer (BC) is a savage disease, executing a huge number of individuals consistently. Creating robotized dangerous BC recognition framework connected on patient's symbolism can assist managing this issue all the more effectively, making diagnosis more versatile and less inclined to mistakes. DeCAF

(or deep) highlights comprise of an in the middle of arrangement it depends on reusing a formerly trained CNN just as highlight vectors, which is then utilized as contribution for a classifier prepared just for the new order assignment. In the light of this, they display an assessment of DeCaf highlights for BC recognition, with a specific end goal to all the more likely see how they contrast with alternate methodologies [1].

This work proposes to classify breast cancer histopathology images independent of their magnifications using convolutional neural networks (CNNs). They propose two different architectures; single task CNN is used to predict malignancy and multi-task CNN is used to predict both malignancy and image magnification level simultaneously. Evaluations and comparisons with previous results are carried out on BreaKHis dataset [2].

The reason for this work is to create an insightful remote discovery and finding approach for breast disease in light of cytological pictures. Initially, this work exhibits a completely mechanized methodology for cell nuclei recognition and division in bosom cytological pictures. The areas of the cell cores in the picture were identified with roundabout Hough change. The expulsion of false-positive (FP) discoveries (loud circles and platelets) was achieve utilizing Otsu's thresholding procedure and fluffy c-implies grouping strategy. The division of the nuclei limits was proficient with the utilization of the marker-controlled watershed change. Next, an astute breast malignancy grouping framework was created [3].

The effectiveness of the treatment of breast cancer depends on its timely detection. An early advance in the finding is the cytological examination of breast material acquired straightforwardly from the tumor. This work gives in PC supported breast growth recognizable proof of issue in light of the examination of cytological pictures of fine needle biopsies to recognize this biopsy as either benevolent or harmful. Rather than confer on the exact division of cell nuclei, the nuclei are finding by circles utilizing the roundabout Hough change system. The result circles are then sifted to keep just astounding estimations for additionally think about by a help vector machine which groups identified circles as right or wrong utilizing surface highlights and the level of cores pixels as per a cores veil acquired utilizing Otsu's thresholding system [4].

This work direct some fundamental examinations utilizing the deep learning way to deal with arrange breast cancer histopathological pictures from BreaKHis, an openly dataset accessible at <http://webinf.ufpr.br/vri/bosom> malignancy database. They propose a strategy in view of the extraction of picture patches for preparing the CNN and the mix of these patches for definite grouping. This strategy means to permit utilizing the high-goals histopathological pictures from BreaKHis as contribution to existing CNN, maintaining a strategic distance from adjustments of the model that can prompt a more unpredictable and computationally exorbitant engineering [5].

Current methodologies depend on handcraft highlight portrayal, for example, shading, surface, and Local Binary Patterns (LBP) in arranging two areas. Contrasted with

carefully assembled include based methodologies, which include undertaking subordinate portrayal, DCNN is a conclusion to-end highlight extractor that might be straightforwardly gained from the crude pixel force estimation of EP and ST tissues in an information driven mold. These abnormal state highlights add to the development of a directed classifier for separating the two kinds of tissues [6].

The test turns out to be the means by which to cleverly join fix level arrangement results and model the way that not all patches will be discriminative. They propose to prepare a choice combination model to total fix level forecasts given by fix level CNNs, which to the best of our insight has not been appeared previously. They apply the technique to the grouping of glioma and non-little cell lung carcinoma cases into subtypes [7].

Computerized atomic identification is a basic advance for various PC helped pathology related picture examination calculations, for example, for mechanized evaluating of breast disease tissue examples. Be that as it may, computerized core location is muddled by (1) the huge number of nuclei and the measure of high goals digitized pathology pictures, and (2) the inconsistency in estimate, shape, appearance, and surface of the individual nuclei. As of late there has been enthusiasm for the utilization of "Profound Learning" techniques for order and investigation of enormous picture information [8].

This work present a dataset of 7,909 breast tumor (BC) histopathology pictures procured on 82 patients, that is currently openly accessible from <http://web.inf.ufpr.br/vri/breast-cancer-database>. The dataset incorporates both benign and malignant pictures. The undertaking related to this dataset is the robotized classification of these pictures in two classes, which would be an important PC helped finding instrument for the clinician. So as to evaluate the trouble of this undertaking, we demonstrate some primer outcomes acquired with state-of-the-art image classification systems [9].

There are a few issues still exist in conventional individual Breast Cancer Diagnosis. To take care of the issues, an individual credit appraisal display in view of help vector order technique is proposed. Utilizing SPSS Clementine information mining device, the individual credit information is bunching investigation by Support Vector Machine. It is investigated in detail with the distinctive part capacities and parameters of Support vector machine. Bolster vector machine could be utilized to enhance crafted by medicinal specialists in the determination of breast growth [10].

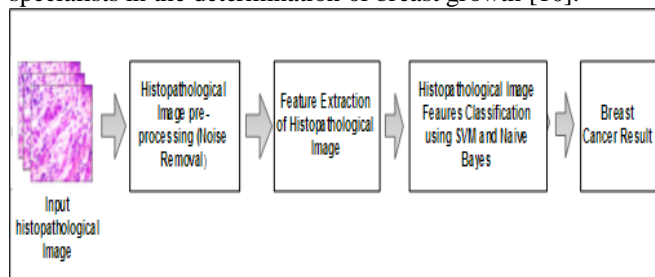


Fig. 1. Proposed System Architecture

III. PROPOSED METHODOLOGY

Classifying breast cancer histopathological images automatically is an important task in computer assisted pathology analysis. However, extracting informative and non-redundant features for histopathological image classification is challenging.

In our proposed work using Histopathological image, firstly we will apply image pre-processing technique to remove the noise of an image. After that we will apply the feature extraction process. The feature-based approaches consist of the features extraction phase and then classification phase. This approach focuses on extracting the feature of image and classify them using machine learning classification technique. The extracted features are trained using support vector machines and Naive Bayes Classification technique. Finally, we compared the performance using the existing classification methods.

Advantages of Proposed System:

1. Work could be beneficial to obtain fast and precise quantification, reduce observer variability, and increase objectivity.
2. Cell nuclei detection using image thresholding and image edge detection.
3. We can measure accurate cell features.
4. This application can be used by physicians from their homes or any other place.
5. This work will be suitable for images with a high degree of noise and blood cells and cell overlapping, as it can successfully detect the cell nuclei.

A. Architecture

Explanation: Input of system:

In this proposed system, we take histopathological breast image as an input for processing.

Image Pre-processing:

In this step, check the size of input image and then the input image is converted into grayscale image. Also, we remove the noise of image using noise reduction technique that i.e. here we use the median filter for noise reduction.

Feature Extraction: CL_{SS_i}

In this step, after image preprocessing, we extract all feature of preprocessed image i.e. infected and healthy cell nuclei.

Classification:

In this step, after image feature extraction, we classify the infected and healthy cell nuclei using support vector machine and also naive bays classification technique.

Result:

This step displays the final breast cancer result.

B. Algorithms

1. Support Vector Machine:

Support Vector Machine (SVM) is used to classify the fruit SS_i $i=1$

Or weighted mean value:

$$\bar{CL} = \sum_{i=1}^k W_i CL_i$$

quality. SVM Support vector machines are mainly two class classifiers, linear or non-linear class boundaries.

The idea behind SVM is to form a hyper plane in between the data sets to express which class it belongs to.

The task is to train the machine with known data and then SVM find the optimal hyper plane which gives maximum distance to the nearest training data points of any class **Steps:**

- Step 1: Read the test image features and trained features.
- Step 2: Check the all test features of image and also get all train features.
- Step 3: Consider the kernel.
- Step 4: Train the SVM using both features and show the output.
- Step 5: Classify an observation using a Trained SVM Classifier.

2. Nave Bays Classification:

Naive Bayes algorithm is the algorithm that learns the probability of an object with certain features belonging to a particular group/class. In short, it is a probabilistic classifier.

The Naive Bayes algorithm is called naive because it makes the assumption that the occurrence of a certain feature is independent of the occurrence of other features.

The Naive Bayesian classifier is based on Bayes theorem with the independence guess between predictors.

A Naive Bayesian model is easy to form, with no critical iterative parameter computation which makes it particularly useful for very large datasets.

Regardless of its simplicity, the Naive Bayesian classifier often does particularly well and is widely used because it often outperforms more experienced classification methods.

C. Mathematical Model

1. Mathematical Equations of Support Vector Machine:

We have k sub-spaces so that there are k classification results of sub-space to classifying breast cancer cells, called $CL_{SS1}, CL_{SS2}, \dots, CL_{SSk}$. Thus the problem is how to integrate all of those results. The simple integrating way is to calculate the mean value:

$$CL = \frac{1}{k} \sum_{i=1}^k CL_{SS_i} \tag{1}$$

$P(A)$ and $P(B)$: Probabilities of the occurrence of event A and B respectively

$P(\frac{B}{A})$: Probability of the occurrence of event B given the event A is true

IV. RESULT AND DISCUSSION

Experiments will be done by a personal computer with a configuration: Intel (R) Core (TM) i5-6700HQ CPU @ 2.60GHz, 16GB memory, Windows 8, MySQL Server 5.1 and Jdk 1.8. Some functions used in the algorithm are provided by Opencv2.4.7.

We will use the histopathological breast image data set

$$\bar{CL} = \sum_{i=1}^k W_i CL_i \tag{2}$$

which consist of total 500 images, patient of collected from the 500 different both healthy

Where W_i is the weight of classification result of subspace, i.e. breast cancer cells result, SS_i and satisfies:

$$\sum_{i=1}^k W_i = 1 \tag{3}$$

The centroid is calculated as follows:

$$X = \frac{\sum_{i=0}^k x_i}{k}, Y = \frac{\sum_{i=0}^k y_i}{k} \tag{4}$$

Where (\bar{X}, \bar{Y}) represents the centroid of the hand, X_i and Y_i are x and y coordinates of the i th pixel in the hand region and k denotes the number of histopathological image pixels that represent only the hand portion.

In the next step, the distance between the centroid and the pixel value was calculated. For distance, the following Euclidean distance was used:

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \tag{5}$$

Where (x_1, x_2) and (y_1, y_2) represent the two co-ordinate values of histopathological image pixel.

2. Mathematical equation in Naive-Bayes Classification:

It gives us a method to calculate the conditional probability, i.e., the probability of an event based on previous knowledge available on the events. Here we will use this technique for breast cancer classification. More formally, Bayes' Theorem is stated as the following equation:

$$P\left(\frac{A}{B}\right) = \frac{P(B)P(A)}{P(A)} \tag{6}$$

Let us understand the statement first and then we will look at the proof of the statement. The components of the above statement are:

$P\left(\frac{A}{B}\right)$: Probability (conditional probability) of occurrence of event A given the event B is true

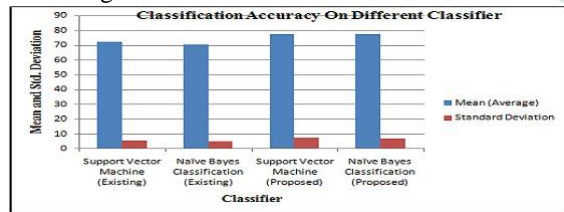


Fig. 2. Classification Accuracy Graph

and infected. This data set collected from city hospital.

The mean and standard deviation values of the input image are computed in each spectral channel as the feature. We let n be the number of pixels in the input image, and v_{ij} denotes the j th band value of the i th pixel in a image. The mean (*mean*) and standard deviation (*std*) of the patch are calculated according to

$$\text{Mean} = \frac{\sum_{i=1}^n v_i}{n} \tag{7}$$

$$\text{Std} = \sqrt{\frac{\sum_{i=1}^n v_i^2 - \frac{(\sum_{i=1}^n v_i)^2}{n}}{n}} \tag{8}$$

Table I is a summary of classification accuracies among different classifiers based on the feature for classifiers. Note that the Support Vector Machine and Nave Bayes-based classifier outperform other classifiers. The classification accuracy for SVM and NB is 77.5% and 77.2% on average, respectively.

TABLE I
TABLE OF CLASSIFICATION ACCURACY

Classification Accuracy				
Classifier	Mean (Exis. System)	Standard Dev. (Exis. System)	Mean (Prop. System)	Standard Dev. (Prop. System)
Support Vector Machine	72.1	5.8	77.5	7.4
Nave Bayes Classification	70.3	5.0	77.2	7.1

V.

CONCLUSION

In this work, we work on histopathological images by using Support Vector Machine (SVM) and Naive Bayes Classification with various configurations for the classification of breast cancer histology images into benign and malignant. The designed SVM topology and Naive Bayes Classification worked well on histopathological images features in classification tasks. However, the performance of the SVM classification and Naive Bayes Classification are better compared to the one of the existing classification methods. SVM have become state-of-the-art, demonstrating an ability to solve challenging classification tasks. This proposed work successfully classifies using breast cancer histology images into benign and malignant.

ACKNOWLEDGMENT

The authors would like to thank the researchers as well as publishers for making their resources available and teachers for their guidance. We are thankful to the authorities of Savitribai Phule University of Pune and concern members of cPGCON2018 conference, organized by, for their constant guidelines and support. We are also thankful to the reviewer for their valuable suggestions. We also thank the college authorities for providing the required infrastructure and support. Finally, we would like to extend a heartfelt gratitude to friends and family members.

REFERENCES

- [1] F. A. Spanhol, L. S. Oliveira, P. R. Cavalin, C. Petitjean, and L. Heutte, Deep features for breast cancer histopathological image classification, 2017 IEEE International Conference
- [2] N. Bayramoglu, J. Kannala, and J. Heikkila, Deep learning for magnification independent breast cancer histopathology image classification, 2016 23rd International Conference on Pattern Recognition (ICPR), 2016.
- [3] Y.M. George, H. H. Zayed, M. I. Roushdy, and B. M. Elbagoury, Remote computer-aided breast cancer detection and diagnosis system based on cytological images, IEEE Systems Journal, vol. 8, no. 3, pp. 949964, Sept 2014.
- [4] P. Filipczuk, T. Fevens, A. Krzyzak, and R. Monczak, Computeraided breast cancer diagnosis based on the analysis of cytological images of fine needle biopsies, IEEE Transactions on Medical Imaging, vol. 32, no. 12, pp. 21692178, 2013.
- [5] F. A. Spanhol, L. S. Oliveira, C. Petitjean, and L. Heutte, Breast cancer histopathological image classification using convolutional neural networks, in International Joint Conference on Neural Networks, Vancouver, BC, Canada, July 2016, pp. 25602567.
- [6] J. Xu, X. Luo, G. Wang, H. Gilmore, and A. Madabhushi, A deep convolutional neural network for segmenting and classifying epithelial and stromal regions in histopathological images, Neurocomputing, vol. 191, pp. 214223, 2016.
- [7] L. Hou, D. Samaras, T. M. Kurc, Y. Gao, J. E. Davis, and J. H. Saltz, Patch-based convolutional neural network for whole slide tissue image classification, in IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, Nevada, USA, June 2016, pp.24242433.
- [8] J. Xu, L. Xiang, Q. Liu, H. Gilmore, J. Wu, J. Tang, and A. Madabhushi, Stacked sparse autoencoder (SSAE) for nuclei detection of breast cancer histopathology images, IEEE transactions on medical imaging, vol. 35, no. 1, pp. 119130, 2016.
- [9] F. Spanhol, L. Oliveira, C. Petitjean, and L. Heutte, A dataset for breast cancer histopathological image classification, IEEE Transactions on Biomedical Engineering, vol. 63, no. 7, pp. 14551462, 2016.
- [10] Shang Gao Hongmei Li, Breast Cancer Diagnosis Based on Support Vector Machine. 2012 International Conference on Uncertainty Reasoning and Knowledge Engineering
- [11] Su H, Liu F, Xie Y, Xing F, Meyyappan S, Yang L. Region segmentation in histopathological breast cancer images using deep convolutional neural network. In Biomedical Imaging (ISBI), 2015 IEEE 12th International Symposium on 2015 Apr 16 (pp. 55-58).
- [12] Wan, S., Huang, X., Lee, H.C., Fujimoto, J.G. and Zhou, C., 2015, April. Spoke-LBP and ring-LBP: New texture features for tissue classification. In Biomedical Imaging (ISBI), 2015 IEEE 12th International Symposium on (pp. 195-199). IEEE.
- [13] Levi, Gil, and Tal Hassner. Age and gender classification using convolutional neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 34-42. 2015.
- [14] M. Veta et al., Breast cancer histopathology image analysis: A review, IEEE Transactions on Biomedical Engineering, vol. 61, no. 5, p. 1400, 1411 2014.
- [15] C. Desir et al., Classification of endomicroscopic images of the lung based on random subwindows and extra-trees, IEEE Transaction on Biomedical Engineering, vol. 59, no. 9, pp. 26772683, 2012.