# A Novel Feature Extraction Technique for Content Based Image Classification in Digital Marketing Platform

Rik Das

Assistant Professor, Dept. of Information Technology
Xavier Institute of Social Service
Ranchi, India
rikdas78@gmail.com

Dr. Subhajit Bhattacharya

Assistant Professor, Dept. of Marketing Management
Xavier Institute of Social Service
Ranchi, India
tobesubho@gmail.com

*Abstract*— **Increasing size and complexity of information about different products offered through digital media has made it essential to target the effective and efficient techniques of data management and categorization. Text data has gradually been replaced by image or visual data due to increasing importance of image capturing devices and social media. Images have been used as one of the foremost linkages between the brand and the consumer in digital marketing domain. Classification has been considered as a vital component of machine learning necessary for data identification which can be initiated before retrieval to restrict the search within the class of interest. The authors have proposed a novel technique for feature extraction for content based image classification. Analyses of customer satisfaction related to content based product identification with images in diverse media have also been carried out. The classification results with the proposed technique have outperformed the existing methods and have shown an increment of 20% in precision results.**

**Keywords—Feature Extraction, Content Based Image Classification, Binarization, Digital Marketing**

## I. INTRODUCTION

Visual data has strong impact on consumer attention that has kindled the buying intention of a consumer [1,2]. Advent of globalization has rapidly influenced the customer preferences and demands [3,4]. Prospective client support has been significantly based on consumer satisfaction [5,6]. Consumer satisfaction has lead to strong and positive behavioural outcomes connected to sustainable purchase [7].

The revenue of the business process has been adversely affected by predominant dissatisfaction of consumers due to large amount of irrelevant results generation for a text based query . The customers were not getting a perishable searching environment. This has led to the popularity of visual data which has principally reinstated the text based keyword searching both online and offline. The authors have proposed a simple but efficient technique of feature extraction from image data and have compared the classification result with proposed technique of feature extraction to that of the existing techniques. Initially, the classification process was carried out with three different classifiers followed by ensemble of classifiers. The results for classification with the novel method have outclassed the prevailing techniques and have significantly enhanced the classification accuracy.

## II. RELATED WORK

Accessibility of online and offline computer based diverse information for products and services has drastically stimulated the customer interactions [19, 20]. Traditional means to locate the product of interest by the customers has been based on text queries. However, the method has huge amount of irrelevant results as output. One of the driving factors for inappropriate output has been due to reprehensible selection of keywords as query. Recent approaches of searching has emphasized upon the content of the searched object rather than its name as a keyword [21, 22]. The content based searching process has been facilitated by the product

image which can provide the necessary knowledge for the required product based on its visual contents and has been anticipated to filter out the unwanted results with higher probability. The content based process of image classification has been carried out by interacting with image data in terms of features which has evidently described the intrinsic image property. Binarization process of differentiating the image foreground from its background has proven efficiency to facilitate image feature extraction. Images have been affected by varied illumination and incoherent gray levels which can adversely influence the binarization process. Hence, the process of binarization depends on selection of appropriate thresholds which may be global threshold [16], local threshold [11,15,17] or mean threshold [8,9,18] for different image qualities. The process of mean threshold selection has the least computation overhead and faster execution compared to the other two techniques. The irregular illumination of stained images was dealt resourcefully during feature extraction by considering measures of dispersion like standard deviation and variance for threshold selection to binarize the images for feature extraction by Thepade et. al, 2014 [10], Ramírez-Ortegón, M.A. and Rojas R., [12], Liu.C [13] and Shaikh [14]. The authors have proposed a novel technique of feature extraction by binarization with mean threshold. The technique has considered the contribution of mean and standard deviation for calculation of feature vectors and has improved the classification performance with small feature vector size.

## III. Research Conduction

The authors have focused on connecting the marketing research domain to technology improvement for a novel feature extraction technique for content based object classification. It has created a new platform of client satisfaction with augmented accuracy in online shopping and digital marketing. The authors have presented a method in algorithmic form aimed at feature extraction for image data classification for enhancing consumer satisfaction related to product recognition. Initially, the statistical analyses have been done to check the consumer satisfaction level connected to product image classification. Image or visual data has been commonly used for current digital platform to locate the product of choice. Classification of image as a precursor of retrieval can reduce irrelevant result generation by restricting the search within the category of interest. An empirical set of 237 sample responses has been collected through e-mail. The sampling was completed using a list-based sampling frame [23, 24]. The Likelihood Ratios, Chi-Square Tests (Phi, Cramer's V tests) had been used to determine the relationship between the consumer dissatisfaction in product identification in digital media pertaining to dissimilar outcomes, large number of options, inaccuracy & inefficiency. IBM SPSS (version 20) has being used for statistical testing and analysis. The authors have proposed a novel method of feature extraction to facilitate image classification which can act as an antecedent for retrieval to enhance consumer satisfaction with product recognition at different media.

## IV. Hypothesis Testing

Satisfaction is considered as one of the most important key components of business success through proper consumer acquisition and management. In this respect consumer opinion analysis is important to determine the factor that could create dissatisfaction in consumer. The authors have primarily identified and tested the contributions of the several components relating to consumer dissatisfaction in connection with product identification. The factors which had been considered are time complexity, dissimilar outcomes, large number of options inaccuracy and inefficacy. Some specified hypothesis have been formulated and tested to support the current research agenda. The authors have proposed a novel technique for feature extraction for image classification which can be used as a precursor of retrieval to increase consumer satisfaction in product recognition at digital marketing platform.

**Hypothesis 1:** *Consumer dissatisfaction in product retrieval at in different Media has been associated with time taken in searching process, dissimilar outcome , large number of options as output and inaccuracy & inefficacy.*

Analysis in Table- 1 has shown significant association in between Consumer dissatisfaction in product retrieval in different medias and time complexity (Likelihood Ratio=59.686; Phi=0.672; Cramer's V=0.388and p<0.01). In the same way dissimilar outcome with searching key word (Likelihood Ratio=61.379; Phi=0.497; Cramer's V= 0.249 and p<0.01) and inaccuracy & inefficacy (Likelihood Ratio=120.065; Phi=0.712; Cramer's V= 0.411and p<0.01) have significant association with Consumer dissatisfaction in product retrieval in different medias. But the large number of appropriate options as output (Likelihood Ratio=5.55; Phi=0.142; Cramer's V= 0.100and p> 0.05) have not adversely concerned the consumer dissatisfaction in product identification in different medias.

understood as the key factor which has caused consumer dissatisfaction in product identification from different media. From From the aforesaid statistical calculation it can be observed that consumers as well as client dissatisfaction in product identification in different medias was caused by time complexity, dissimilar and irrelevant outcomes and inaccuracy

| Chi-Square Tests | | | | |
|---|---|---|---|---|
| **Raw variable → Consumer dissatisfaction in product classification in different medias** | | | | |
| **Column variables ↓** | *Likelihood Ratio* | *Sig. (2-sided)* | *Phi* | *Cramer's V* |
| Time taken in searching process | 59.686 | .000 | 0.672 | 0.388 |
| Dissimilar outcome   with searching key word | 61.379 | .000 | 0.497 | 0.249 |
| Large number of options as output | 5.55 | 0.781 | 0.142 | 0.100 |
| Inaccuracy & inefficacy | 120.065 | .000 | 0.712 | 0.411 |
| *Number of sample 237* | | | | |

Table1. Chi-Square Test for Consumer dissatisfaction analysis

| **Kendall's tau: non-parametric correlation coefficient** Table | | Consumer dissatisfaction in product classification in different medias |
|---|---|---|
| Time taken in searching process | Correlation Coefficient | .298** |
| | Sig. (2-tailed) | .000 |
| Dissimilar outcome        with searching key word | Correlation Coefficient | .345** |
| | Sig. (2-tailed) | .000 |
| Large number of options as output | Correlation Coefficient | -0.004 |
| | Sig. (2-tailed) | 0.664 |
| Inaccuracy & inefficacy | Correlation Coefficient | .530** |
| | Sig. (2-tailed) | .000 |
| *\*\*. Correlation is significant at the 0.01 level (2-tailed); Number of sample 237* | | |

Table2. Nonparametric Correlation Test for Consumer Dissatisfaction Analysis

From the above results shown in Table 2, it was observed that consumer dissatisfaction in product identification in different media was having a strong association with the time taken in searching process and dissimilar outcome for searching by keywords. On the other hand the consumer dissatisfaction in product retrieval was not associated with large number of appropriate options as output. Inaccuracy & inefficacy can be

and inefficacy. The proposed method has stimulated optimized outcome for product recognition by implementing a novel feature selection method for object recognition and has amalgamated the concept in the field of marketing related to online or offline product recognition. The aforesaid statistical calculation has revealed that consumers as well as client dissatisfaction in product identification in different medias

was caused by time complexity, dissimilar and irrelevant outcomes and inaccuracy & inefficacy. The proposed method has stimulated optimized outcome for product recognition by implementing a novel feature selection method for object recognition and has amalgamated the concept in the field of marketing related to online or offline product recognition.

**Hypothesis 2:** *Visual data based query and text data based query in computer based digital media platform have similar impact on consumer satisfaction connecting to product recognition.*

From the above results shown in Table 4, it was observed that consumer satisfaction and Visual data based query related to product recognition in diversified product categories in computer based digital media platforms are significantly correlated.

On the other hand it was detected that insignificant correlation lies in between consumer satisfaction and Text data based query related to product recognition in diversified product categories in computer based digital media platforms.

**Chi-Square Tests**

**Raw variable →** Consumer satisfaction related to product recognition in diversified product categories in computer based digital media platforms

| **Column variables ↓** | *Likelihood Ratio* | *Sig. (2-sided)* | *Phi* | *Cramer's V* |
|---|---|---|---|---|
| Visual data based query | 449.854 | .000 | 2.00 | 1.00 |
| Text data based query | 13.711 | .090 | 0.241 | 0.170 |
| *Number of sample 237* | | | | |

Table3. Chi-Square Test for Consumer Satisfaction Analysis

| **Kendall's tau: non-parametric correlation coefficient** | | Consumer satisfaction related to product recognition in diversified product categories in computer based digital media platforms |
|---|---|---|
| Visual data based query | Correlation Coefficient | .993** |
| | Sig. (2-tailed) | .000 |
| Text data based query | Correlation Coefficient | .009 |
| | Sig. (2-tailed) | .877 |
| ***. Correlation is significant at the 0.01 level (2-tailed); Number of sample 237* | | |

Table4. Nonparametric Correlation Test for Consumer Dissatisfaction Analysis

Analysis in Table-3 has shown significant association among consumer satisfaction and Visual data based query (Likelihood Ratio=449.854; Phi=2.00; Cramer's V=1.00 and $p<0.01$) .Related to product recognition in diversified product categories in computer based digital media platform with text data based query is not having significant relationship (Likelihood Ratio=13.711; Phi=0.241; Cramer's V= 0.170 and $p> 0.05$) . From the analysis it was observed that visual data based query is having better association with consumer satisfaction related to product recognition compare to Text data based query.

## V. PROPOSED TCHNIQUE

The proposed technique of feature extraction was initiated with the extraction of Red (R), Green (G) and Blue (B) color components from the images. Individual mean threshold values were calculated for each of the color component as shown in equation 1.

$$Tavx = (1/m*n)*\sum_{i=1}^{m}\sum_{j=1}^{n} x(i,j) \tag{1}$$

x= R, G and B respectively for each of the corresponding color component considered

Threshold selection process was followed by calculation of binary bitmaps as shown in equation 2. The process of bitmap selection has assigned a value 1 to the pixel values higher than or equal to the threshold and has allocated a value 0 in case the pixel value is lesser than the threshold.

$$BitMap_x = \begin{cases} 1, & iff \dots x(i,j) >= Tav_x \\ 0 & iff \dots x(i,j) < Tav_x \end{cases} \quad (2)$$

Two feature vectors namely the higher intensity feature vector and the lower intensity feature vector was computed from the pixel values assigned with 1 and 0 respectively for each color component. The process of determining the feature vector has been given in equation 3 and 8 where the mean and

the standard deviation of each of the clusters of pixels were derive the final signature from the images.

$$xhi_{mean} = mean \sum_i \sum_j (x(i,j)) > Tav_x \quad (3)$$

$$xhi_{stdev} = \sigma \sum_i \sum_j (x(i,j)) > Tav_x \quad (4)$$

$$xhi_{F.V.} = xhi_{mean} + \left( xhi_{mean} + xhi_{stdev} \right) \quad (5)$$

$$xlo_{mean} = mean \sum_i \sum_j (x(i,j)) < Tav_x \quad (6)$$

$$xlo_{stdev} = \sigma \sum_i \sum_j (x(i,j)) < Tav_x \quad (7)$$

$$xlo_{F.V.} = xlo_{mean} + \left( xlo_{mean} + xlo_{stdev} \right) \quad (8)$$

Where, x represents R, G and B for individual components and $Tx$ is the threshold value for each pixel.

## VI. EXPERIMENTAL VERIFICATION

A widely used public dataset namely the Wang dataset [18] was considered for the assessment purpose as shown in Fig 1.



Fig. 1 Sample Wang Dataset

calculated to

The process of 10 fold cross validation was carried out in which 9 subsets were considered as training set and 1 subset was considered as the testing set. The final results were asserted averaging the results of the 10 iterations for 10 fold cross validation

Three different classifiers namely K Nearest Neighbor (KNN), Support Vector Machine (SVM) and Artificial Neural

Network (ANN) were used for evaluation purpose. KNN has considered similarity functions of two different instances for classification results. SVM has conducted the learning process of Self Organizing Map (SOM) for classification and has presumed that only nearby nodes has affected the behavior of

each other. Finally, ANN was implemented with a feed forward architecture known as multi layer perceptron (MLP). Individual performance of the classifiers with proposed feature extraction technique was measured by the evaluation metrics called Precision and Recall. Further, the classifiers were ensemble by means of maximum probability and the precision

and recall values were calculated for classification with the proposed method of feature extraction. The comparative evaluation for precision and recall with individual classifiers and the ensemble of classifiers has been illustrated in Table 5 and 6 and in Fig. 3

| Categories | KNN | SVM | ANN | Ensemble |
|---|---|---|---|---|
| Tribals | 85.2 | 78 | 80.2 | 83.3 |
| Sea Beach | 81.9 | 76.5 | 82.5 | 84 |
| Gothic Structure | 60.4 | 63.4 | 73.9 | 73.9 |
| Bus | 59.3 | 67.3 | 70.9 | 70.1 |
| Dinosaur | 100 | 100 | 97.1 | 98 |
| Elephant | 68.9 | 73.2 | 81.2 | 80.7 |
| Roses | 94.2 | 85.7 | 93.9 | 93.9 |
| Horses | 93.8 | 98.9 | 94.9 | 95.9 |
| Mountains | 77.2 | 82.7 | 95.7 | 96.7 |
| Food | 63.7 | 70.9 | 93.7 | 91.8 |
| Average | 78.5 | 79.7 | 86.4 | 86.8 |

Table5. Comparison of Precision

| Categories | KNN | SVM | ANN | Ensemble |
|---|---|---|---|---|
| Tribals | 52 | 64 | 81 | 80 |
| Sea Beach | 68 | 75 | 80 | 79 |
| Gothic Structure | 55 | 52 | 65 | 65 |
| Bus | 73 | 72 | 78 | 82 |
| Dinosaur | 100 | 100 | 100 | 100 |
| Elephant | 91 | 93 | 95 | 96 |
| Roses | 81 | 90 | 93 | 93 |
| Horses | 91 | 90 | 93 | 94 |
| Mountains | 71 | 81 | 88 | 87 |
| Food | 86 | 78 | 89 | 89 |
| Average | 76.8 | 79.5 | 86.2 | 86.5 |

Table6. Comparison of Recall

**Comparison of Precision and Recall for Classification with Proposed Feature Extraction Technique for different classifier environments**



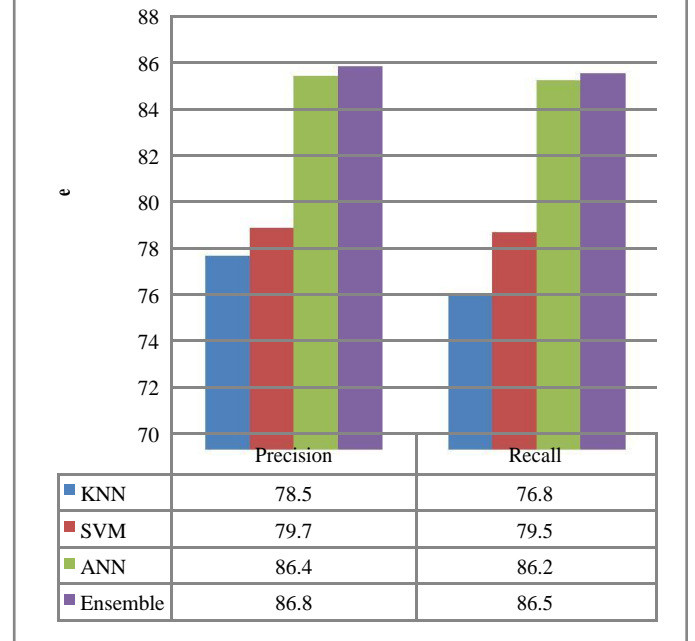| | Precision | Recall |
|---|---|---|
| KNN | 78.5 | 76.8 |
| SVM | 79.7 | 79.5 |
| ANN | 86.4 | 86.2 |
| Ensemble | 86.8 | 86.5 |

Fig. 2 Comparison of Precision and Recall

Henceforth, the results of classification with the proposed feature extraction technique were compared to the state-of-the art techniques.

The illustration in Fig. 2 has clearly revealed the supremacy of proposed technique of feature vector extraction for classification performances with respect to the existing techniques. Moreover, it was observed from the results in Table 5 and 6 that the classification performances of the proposed feature extraction technique with individual classifiers were also higher than state of the art techniques as illustrated in Fig. 3. Hence, it was inferred that the proposed method of feature extraction has manifested better classification results than the previous techniques.
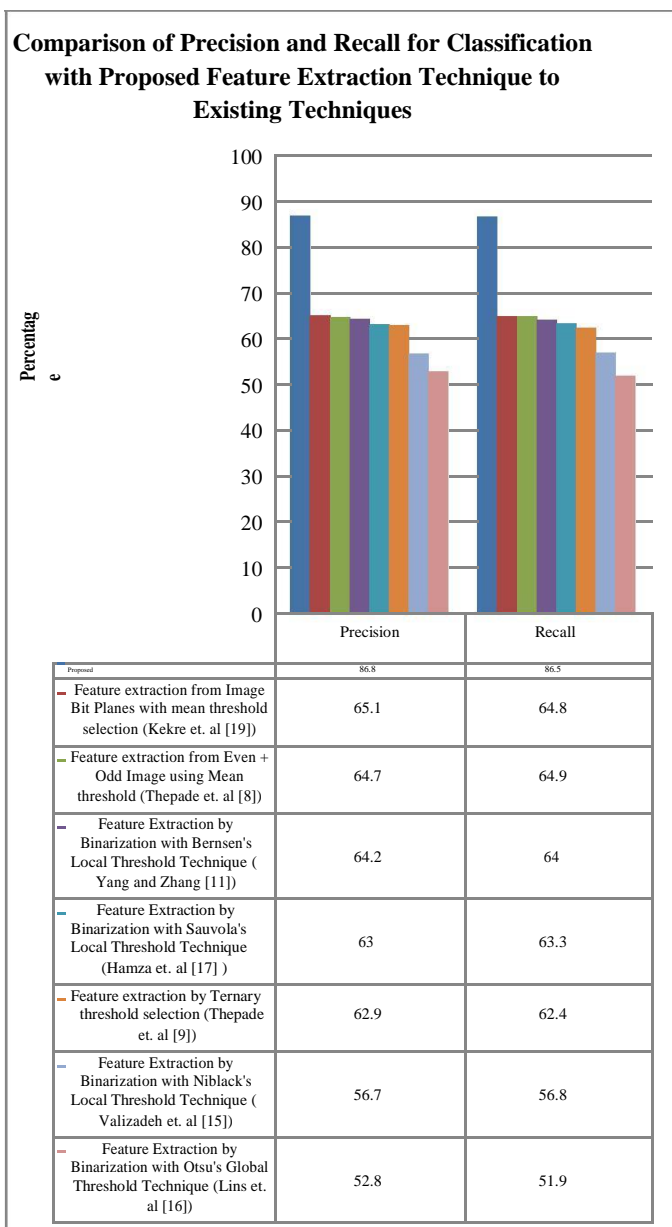
**Comparison of Precision and Recall for Classification with Proposed Feature Extraction Technique to Existing Techniques**

Percentage



| | Precision | Recall |
|---|---|---|
| Proposed | 86.8 | 86.5 |
| Feature extraction from Image Bit Planes with mean threshold selection (Kekre et. al [19]) | 65.1 | 64.8 |
| Feature extraction from Even + Odd Image using Mean threshold (Thepade et. al [8]) | 64.7 | 64.9 |
| Feature Extraction by Binarization with Bernsen's Local Threshold Technique ( Yang and Zhang [11]) | 64.2 | 64 |
| Feature Extraction by Binarization with Sauvola's Local Threshold Technique (Hamza et. al [17] ) | 63 | 63.3 |
| Feature extraction by Ternary threshold selection (Thepade et. al [9]) | 62.9 | 62.4 |
| Feature Extraction by Binarization with Niblack's Local Threshold Technique ( Valizadeh et. al [15]) | 56.7 | 56.8 |
| Feature Extraction by Binarization with Otsu's Global Threshold Technique (Lins et. al [16]) | 52.8 | 51.9 |

Fig. 3 Comparison of Precision

VII. CONCLUSION

Increasing importance of visual data analytics has made the process of image recognition imperative. Efficient feature extraction has been considered as the crucial factor to define the success rate for image data identification. The authors have proposed a novel technique of feature extraction in this paper to boost up the process of image classification for enhance customer satisfaction in digital marketing. The method has shown higher efficiency in classification compared to the prevailing methods of feature extraction and has contributed significantly to the analysis of image data necessary for enhance customer satisfaction in content based product searching for digital marketing platform. The work can be extended towards image data analysis in all the noteworthy fields like pattern recognition, content based image retrieval, security, media and journalism etc.

## *References*

[1] Alexandra , G. 2005. "A user-friendly marketing decision support system for the product line design using evolutionary algorithms", Decision Support Systems, Vol. 38 No. 4, pp. 495 - 509.

[2] Malhotra, N. 2004. Marketing Research: An Applied Orientation, 4th Ed. Upper Saddle River, N.J.: Pearson Education, Inc.

[3] Smith, S.M. and Albaum, G. S. 2010 "Introduction to Marketing Research," Qualtrics Labs, Inc. (Online Edition).

[4] Smith, S.M., Smith, J. and Allred, C.R. 2006. "Advanced Techniques and Technologies in Online Research." In R. Grover and M. Vriens *Eds.), The Handbook of Marketing Research: uses, Misuses and Future Advances . 132-158. Thousand Oaks, CA: Sage Publications.

[5] Lee J.-E., Jin R., Jain, A. K. and Tong ,W. 2012 "Image retrieval in forensics: tattoo image database application,"IEEE Multimedia,vol.19,no.1,pp.40–49.

[6] Ovidiu DOBRICAN, 2009. Multimedia and Decision-Making Process .Informatica Economica vol. 13. no. 3.

[7] Wamg, X.-F. 2009. "The analyse of development of network shopping in china," Economic Research Guide, vol.40, pp.174–175.

[8] Thepade, S, Das, Rik , Ghosh, S(2013). Performance Comparison of Feature Vector Extraction Techniques in RGB Color Space using Block Truncation Coding or Content Based Image Classification with Discrete Classifiers. 2013 India Conference (INDICON), 2013 Annual IEEE Digital Object Identifier: 10.1109/INDCON.2013.6726053 Publication Year: 2013 , p. 1 – 6

[9] Thepade,S., Das, R., Ghosh,S. (2013). Image classification using advanced block truncation coding with ternary image maps. 2013 Springer International Conference on Advances in Computing, Communication and Control, Communications in Computer and Information Science Volume 361, 2013, p.500-509

[10] Thepade Sudeep(Dr.), Kumar Das Rik Kamal, Ghosh Saurav(2014). A Novel Feature Extraction Technique using Binarization of Bit Planes for Content Based Image Classification (In Press). Journal of Engineering, Hindawi Publishing Corporation

[11] Yanli Y. and Zhenxing Z. (2012). A novel local threshold binarization method for QR image, IET International Conference on Automatic Control and Artificial Intelligence (ACAI 2012), p. 224-227

[12] Ramírez-Ortegón, M.A. And Rojas R. (2010). Unsupervised Evaluation Methods Based on Local Gray-Intensity Variances for Binarization of Historical Documents, IEEE 20t. International Conference on Pattern Recognition (ICPR), p. 2029-2032

[13] Liu.C (2013). A new finger vein feature extraction algorithm, IEEE 6th. International Congress on Image and Signal Processing (CISP), 1, p. 395-399

[14] Shaikh, S. H., Maiti, A. K., & Chaki, N. (2013). A new image binarization method using iterative partitioning. Machine Vision and Applications, 24(2), p. 337-350.

[15] Valizadeh, M., Armanfard, N., Komeili, M., Kabir E. (2009): A novel hybrid algorithm for binarization of badly illuminated document images. In: 14th International CSI Computer Conference (CSICC), p. 121–126

[16] Lins, R. D., Simske, S. J., Fan, J., Sá, P., Silva, G. P., Shaw, M., & Thielo, M. (2009). Image classification to improve printing quality of mixed-type documents. Proceedings of the International Conference on Document Analysis and Recognition, ICDAR (pp. 1106-1110)

[17] Hamza, H., Smigiel, E., & Belaid, A. (2005). Neural based binarization techniques. Proceedings of the International Conference on Document Analysis and Recognition, ICDAR (Vol. 2005, pp. 317-321)

[18] Kekre H.B., Thepade S, Das R, Ghosh S (2012). Performance Boost of Block Truncation Coding based Image Classification using Bit Plane

Slicing . International Journal of Computer Applications 47(15):45-48, June 2012

[19] Arabie P, Hubert L .1995. Advances in cluster analysis relevant to marketing research. Stud Classif Data Anal Knowl Org 6:3–19

[20]  Baier D, Gaul W .1999. Optimal product positioning based on paired comparison data. J Econ 89(1):365–392

[21]  Gaul W, Baier D .1994. Marktforschung und Marketing Management: Computerbasierte Entscheidungsunterstützung. Oldenbourg, München

[22] Punj G, Stewart DW .1983. Cluster analysis in marketing research: review and suggestions for application. Journal of Marketing Research 20:134–148

[23]  Couper, Mick P. (2000) 'Review: Web Surveys: A Review of Issuesand Approaches', The Public Opinion Quarterly, 64(4): 464-494.

[24] Fahmy, S., Fosdick, S.B., and Johnson, T.J. (2005) 'Is Seeing Believing? A Survey of Magazine Professionals' Practices andAttitudes Toward EthicalStandardsfor Photographs', Journal of        Magazine andNew Media Research, Spring issue, 1-16.