

A Novel Approach of Character Recognition using Local Binary Pattern and SVM Classifier

Alka Kumari
akalkakumari25@gmail.com

Nirupama Tiwari
girishniru@gmail.com

Abstract –The purpose of this paper is to initiate the development of a framework for digitization and recognition of handwritten character. The software architecture of the system is based on the generic integration of commercial character recognition tools called OCR (Optical Character Recognition). Indeed, the evolution of these tools today opens up a serious alternative to manual entry in order to increase productivity. We will present the capabilities of these tools and their limitations. In this paper, support vector machine (SVM) based approach is proposed to solve this type of problem which is capable of recognizing character with the help of Local Binary Pattern feature extraction and support vector machine (SVM) based classifier.

Keywords –DMS, LBP, OCR, SVM.

I. INTRODUCTION

The digitization of documents is an important step in the implementation of an electronic document management system (DMS). The choice of the scanning solution must take into account all the stages of document processing from the acquisition, the conversion of the content to the correction and the exploitation of the final document. The purpose of this scan is to use the converted content, for example to search for information [1].

The central part of the digitization concerns the character recognition and the structuring of the content. It can be done in two ways, one not exclusive to the other [2].

It is now possible to have this step done manually by operators, at a reduced cost. In fact, more and more people in developing countries are finding a skilled workforce that can do this job perfectly. However, it can be difficult to manage, both for reasons of distance, and because it takes people who can read English. Nor should we neglect the confidentiality issues that may arise [3].

It can be done automatically with tools. It is this aspect that we will develop. Note, however, that you can never permanently remove the user, because it will always require a step of validation and verification of the result provided. The tools (no

more than the human being) will never be perfect, and we tolerate less error on the part of a machine than on the part of a human being. On the other hand, one can expect a significant productivity gain compared to a purely manual input: in fact, the computers can work 24/24, and be bought in sufficient number. This combination operator / machine can only be beneficial, especially since the errors made by a machine (so without interpretation of the content) will not be of the same type as those made by an operator [4-6].

Automatic character recognition techniques have evolved and matured in the last decade and we see the flourishing of the software market increasingly cheaper, more and more complete and reliable offering very reliable retro-conversion solutions. These software packages, belonging to the OCR (Optical Character Recognition) family, are today able to distinguish the different media in the document (text, graphic and photograph), to identify linear and tabular structures, to face a significant variation in typography, interpreting and restoring several editorial styles. Test benches are commonly carried out on the latest versions of these software packages showing the new possibilities offered and giving a correct and accurate assessment of their recognition capacity in terms of confidence rate, accuracy and speed of execution, by type [7].

These tools have also progressed in terms of integration since they exist today in the form of easily integrable APIs in a system or in the form of development kits offering opportunities for expansion and cooperation with other APIs to increase the efficiency of the whole [8].

If these techniques are used today without much difficulty in several fields, their integration into a real digitization chain, in an industrial environment, requires taking into account other constraints related to production such as volumetry, the rate of digitization and above all the high quality of results required at the end of the chain. Few experiments have been reported to give an account of the efficiency of these integrations, but it is clear from

the first tests that we carried out in company, that their contribution is considerable in terms of productivity gain, compared for example to double manual entry.

One of the main problems in choosing a method is knowing how to evaluate the result, according to its needs. In other words, the expected safety of the recognition, or the tolerated error rate is dependent on the intended application. Two examples will allow us to illustrate this fact:

In absolute terms, an error rate of 90% is correct. However, suppose that the application concerns the digitization of the Directory. A recognition rate of 90% is catastrophic because it means an error every 10 characters; a phone number with ten characters, all the numbers will be wrong, and so the recognition will be useless. Currently we are moving towards a mixed recognition, combining human intervention, both for the preparation of documents, verification, correction and complete documents (for example: formulas) and digital processing [9-11].

We propose an optimization method in the context of this work for the selection and weighting of features of a system of recognition of isolated handwritten character. The classification method used is based on the theory of support vector machine (SVM). In this work, the optimization of the recognition system is advanced. The first criterion is to extract the features using Local Binary Pattern, while the second criterion should improve or maintain the performance of the recognition system.

II. PROPOSED METHOD

The flow diagram of proposed character recognition system is shown by Figure 1. The processes are explained below:

A. Data Collection

Scanned image of alphabet data in .JPEG format is loaded to the system. The system we develop here is capable of identifying the character. So the user has to select the character of its own choice for recognition input. Once the input character is selected, several processes are applied to it.

B. Pre-Processing

The goal of the pre-processing is to facilitate the characterization of the form (vowels and consonants) or entity to recognize either by cleaning the picture of the form or reducing the amount of information to process only keep only the most relevant information. The image cleaning is basically to eliminate residual noise from the binarization. Reducing the amount of information to be processed can be obtained from operations to bring the line thickness to a single pixel or by

monitoring feature or from extractors of upper contours, lower and / or interiors.

Note that some forms (vowels and consonants) are inclined or bent so it is necessary to normalize slope this form to segment the form (e.g. segmentation of a word in letters). This standardization is to correct the slope of a word or correct the inclination of the letters in a word to facilitate segmentation.

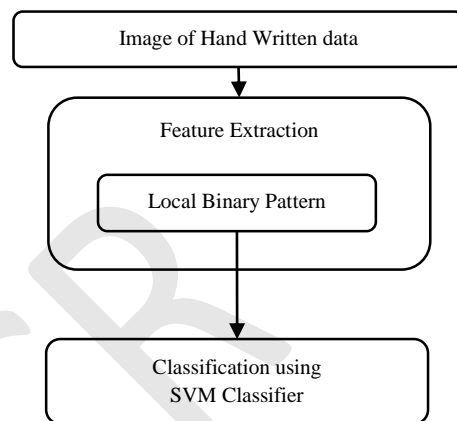


Figure 1: Flow diagram of proposed work

C. Feature Extraction using Local Binary Pattern

Given a central pixel in the image, a pattern code is computed by comparing it with its neighbours:

$$LBP_{P,R} = \sum_{p=0}^{P-1} S(g_p - g_c) 2^p \quad (1)$$

$$S(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (2)$$

Where g_c is the gray value of the central pixel, g_p is the value of its neighbours, P is the total number of involved neighbours and R is the radius of the neighborhood. Suppose the coordinate of g_c is $(0, 0)$, then the coordinates of g_p are $(R * \cos(2\pi p/P), R * \sin(2\pi p/P))$. Figure 2 gives examples of circularly symmetric neighbour sets for different configurations of (P, R) . The gray values of neighbours that are not in the center of grids can be estimated by interpolation.

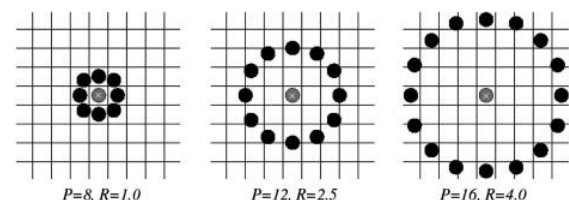


Figure 2: Circularly symmetric neighbour sets for different (P, R)

Suppose the texture image is of size $N \times M$. After identifying the LBP pattern of each pixel (i, j) , a

histogram is built to represent the whole texture image:

$$H(k) = \sum_{i=1}^N \sum_{j=1}^M f(LBP_{P,R}(i,j), k), k \in [0, K] \quad (3)$$

$$f(x, y) = \begin{cases} 1, & x = y \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

Where K is the maximal LBP pattern value. The U value of an LBP pattern is defined as the number of spatial transitions (bitwise 0/1 changes) in that pattern:

$$U(LBP_{P,R}) = |s(g_{P-1} - g_c) - s(g_0 - g_c)| + \sum_{p=1}^{P-1} |s(g_p - g_c) - s(g_{p-1} - g_c)| \quad (5)$$

For example, the LBP pattern 00000000 has a U value of 0 and 01000000 has a U value of 2. The uniform LBP patterns refer to the patterns which have limited transition or discontinuities ($U \leq 2$) in the circular binary presentation. It was verified that only those "uniform" patterns are fundamental patterns of local image texture. In practice, the mapping from $LBP_{P,R}$ to $LBP_{P,R}^{u2}$ (superscript "u2" means that the uniform patterns have a U value of at most 2), which has $P * (P - 1) + 3$ distinct output values, is implemented with a lookup table of 2^P elements. The dissimilarity of sample and model histograms is a test of goodness-of-fit, which could be measured with a nonparametric statistic test.

These feature vectors are classified using the support vector machine classifier.

D. Classification using Support Vector Machine

The classification is developing a decision rule that transforms attributes characterizing the forms in class membership (transition from code space to space-making) [11]. Before a decision model is integrated in a handwriting recognition system, you must have also previously two steps: the learning step and the test step.

As part of our project, support vector machine (SVM) is the method of classification of handwritten character recognition system which is described in the following heading.

Consider the training set $\{x_1, y_1\}, \dots, \{x_\ell, y_\ell\}$, where $x \in X$ and $y \in \{-1, 1\}$, where ℓ is the number of observations and X is a distribution in space \mathcal{R}^n . In the classification problem, the goal is to find an efficient method to construct the optimal separator hyperplane, i.e., with the greatest margin. To do this, one must find the vector w and the constant b , which minimize the norm $|w|^2 = w^T w$ (since it is inversely proportional to the margin), under the constraints:

$$w^T x_i + b \geq 1, \quad \text{if } y_i = 1 \quad (6)$$

$$w^T x_i + b \leq -1, \quad \text{if } y_i = -1 \quad (7)$$

Because one can accept some errors, one relaxes the constraints (6) & (7) and introduces an additional cost related to this relaxation, so that one arrives at the quadratic problem, QP, following:

$$\text{Minimize } \frac{1}{2}(w^T w) + C[\sum_{i=1}^{\ell} \xi_i]$$

w

$$\text{Under the constraints } \begin{cases} y_i(w_i^T x + b) \geq 1 - \xi_i, \\ \xi_i \geq 0 \end{cases} \quad i = 1, \dots, \ell \quad (8)$$

The problem (8) can be solved in the primal space (the space of parameters w and b). In fact, one solves the QP in the dual space, equation (9), (the Lagrange multiplier space) for two main reasons: 1) The constraints (7) and (8) are replaced by the associated Lagrange multipliers, and 2) We obtain a formulation of the problem where the training data appear as an internal product between vectors, which can then be replaced by kernel functions, then construct the hyperplane in the feature space and obtain functions Non-linear in the input space.

$$\text{Maximize } L_D(\alpha) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j (x_i^T x_j) \\ \alpha$$

$$\text{Under the constraints } \begin{cases} \sum_{i=1}^{\ell} y_i \alpha_i = 0, \\ 0 \leq \alpha_i \leq C \end{cases} \quad i = 1, \dots, \ell \quad (9)$$

Where, α_i is the Lagrange multiplier, associated with constraints. Parameter C controls the level of error in the classification.

The SVM evaluation function is defined as:

$$f(x) = \sum_{i=1}^{\ell} \alpha_i y_i k(x_i, x) + b \quad (10)$$

The examples x_i associated with the Lagrange multipliers α_i larger than zero correspond to the support vectors, and have a significant contribution to equation (10). Geometrically, these vectors reside in the margin defined by the separating hyperplane. The constant b represents the threshold of the hyperplane learned in the characteristic space. It can be calculated by the mean of the function (Equation 10), evaluated using the support vectors.

III. SIMULATION AND RESULTS

The performance of proposed algorithms has been studied by means of MATLAB simulation.



Figure 3: Input image

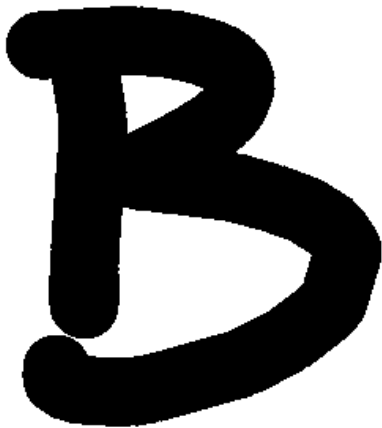


Figure 4: Gray image

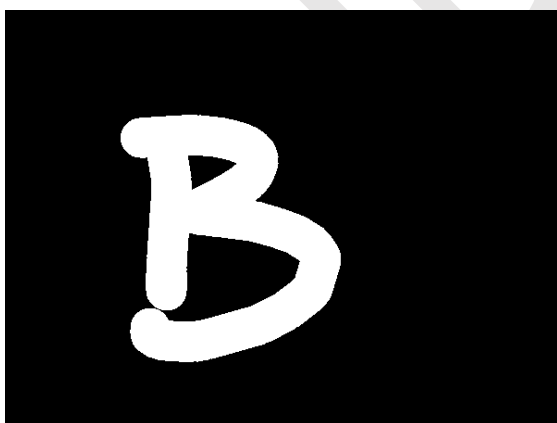


Figure 5: Binary image



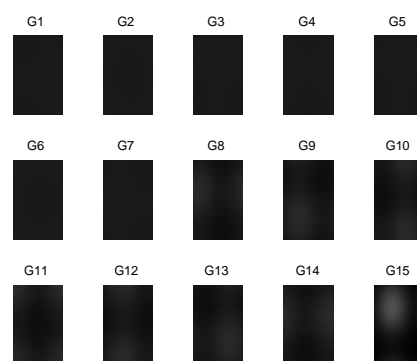
Figure 6: Cropped binary image



Figure 7: Resized binary image



Figure 8: FFT of binary image



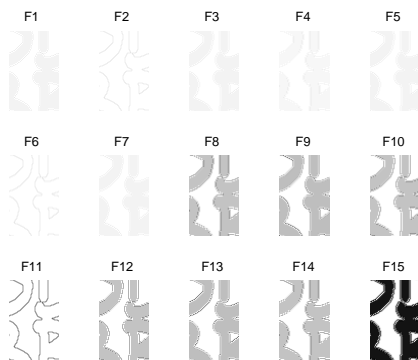


Figure 9: Local binary patterns

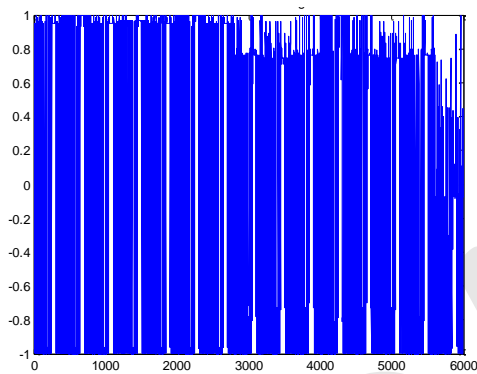


Figure 10: Feature vector of image

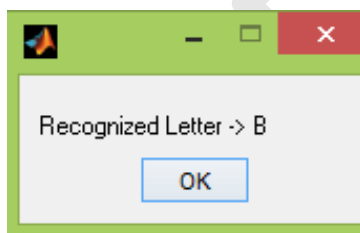


Figure 11: Recognized letter

IV. CONCLUSION

The work presented in this dissertation address the steps necessary to build a system of character recognition. For each of these stages are: pre-processing, extraction of features and classification, have tried to propose an optimization method for the selection of features relevant recognition system. Thus, from the extraction step of the feature, the selection of relevant and non-redundant system is performed features. This selection is to reduce inputs of the classifier (SVM) while improving or maintaining the classification recognition rate. Main contribution of this dissertation is the study of the local binary pattern features. The results of this

selection is satisfactory. The overall accuracy achieved by simulation is 92.6923%

REFERENCE

- [1] Gauri Katiyar, Shabana Mehruz, "MLPNN Based Handwritten Character Recognition Using Combined Feature Extraction", IEEE, International Conference on Computing, Communication and Automation (ICCCA2015), pp. 1155-1159, July 2015.
- [2] Duda, R.O., Hart, P.E. and Stork, D.G., 2012. Pattern classification. John Wiley & Sons.
- [3] Jayadevan, R., Kolhe, S.R., Patil, P.M. and Pal, U., 2011. Offline recognition of Devanagari script: A survey. IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews), 41(6), pp.782-796.
- [4] Kolman, E. and Margalioth, M., 2008. A new approach to knowledge-based design of recurrent neural networks. IEEE Transactions on Neural Networks, 19(8), pp.1389-1401.
- [5] Gaur, A. and Yadav, S., 2015, January. Handwritten Hindi character recognition using k-means clustering and SVM. In Emerging Trends and Technologies in Libraries and Information Services (ETLLIS), 2015 4th International Symposium on (pp. 65-70). IEEE.
- [6] Sameeksha Barve, "Artificial Neural Network Based On Optical Character Recognition", International Journal of Engineering Research & Technology (IJERT), ISSN: 2278-0181, Vol. 1, Issue 4, June 2012.
- [7] Sameeksha Barve, "Optical Character Recognition Using Artificial Neural Network", International Journal of Advanced Technology & Engineering Research (IJATER), ISSN NO: 2250-3536 Volume 2, Issue 2, May 2012.
- [8] Shabana Mehruz, Gauri katiyar, "Intelligent Systems for Off-Line Handwritten Character Recognition: A Review", International Journal of Emerging Technology and Advanced Engineering, ISSN 2250-2459, Volume 2, Issue 4, April 2012.
- [9] Dileep Kumar Patel, Tanmoy Som, Sushil Kumar Yadav, Manoj Kumar Singh, "Handwritten Character Recognition Using Multiresolution Technique and Euclidean Distance Metric", Journal of Signal and Information Processing, PP. 208-214., 2012.
- [10] Rakesh Kumar Mandal, N. R. Manna, "Hand Written English Character Recognition using Column-wise Segmentation of Image Matrix (CSIM)", WSEAS Transactions On Computers, E-ISSN: 2224-2872, Issue 5, Volume 11, May 2012.
- [11] Richard O. Duda, Peter E. Hart, & David G. Stork., "Pattern Classification", (Second edition). New York: Wiley-Interscience, 2001.