

HANDWRITTEN CHARACTER EXTRACTION IN MEDICAL CASE FORMS

Mr. Balaji Murugan P
Department Of Computer Science and
Engineering
SRM Institute Of Science and
Technology
Chennai, India
pbalaji2001@gmail.com

Ms. Manisha E
Department Of Computer Science and
Engineering
SRM Institute Of Science and
Technology
Chennai, India
manopleasancia08@gmail.com

Mr. Vinoth.N.A. S
Department of Computing Technologies
SRM Institute of Science and Technology
Chennai, India
vinoth.nas89@gmail.com

Abstract— Handwritten Character recognition refers to the process of classifying individual characters. It is the ability of an electronic device to read and interpret handwritten input from a variety of sources, including paper documents, digital touch screen devices, and pictures.. This Extraction of handwritten characters from multiple case forms, by checking and dumping them for multiple case forms in the excel sheet for the records is really a tedious process. So to solve this issue , I delve into this problem statement which will drastically resolve the time and cost consumption of the organization. This idea is utilized in a variety of contexts, including form data entry, bank check processing, parcel posting, and mail sorting. However, it is currently a significant problem in the pattern recognition field that will be difficult to fix. Since deep learning is a key methodology in solving detection and pattern recognition problems, various algorithms are available to classify the characters with better prediction rates on various datasets. Ultimately, whichever algorithm produces the best results will be considered the best solution of the character recognition problem. Therefore, various solutions provided by the existing researchers are being discussed using deep learning algorithms in this article. To achieve this, three models namely Keras OCR, Pytesseract, Easter 2.0. were used . Where, Keras and Easter are used to extract the handwritten characters and Pytesseract is used to extract the printed text from the documents or case forms. Finally 76% accuracy is achieved to extract the handwritten characters in the medical forms.

KEYWORDS: Pattern Recognition, Character Extraction and Deep Learning.

I. INTRODUCTION

Now a days in different sectors pattern character recognition is used specially if it is being talked about the handwritten character recognition in posting the parcel, mailed messages sorting, bank related query processing and data form entry. There are various ways of writing the characters and it totally depends on the handwriting style of individuals some of them are making the characters of different size and shape so in the earlier mention application it creates the error to correctly classify the characters due to this it affects the whole process so this is why handwritten recognition is becoming the very notable issue now a days in pattern recognition area.

So, to overcome this issue the main motive is to make handwritten recognition of characters error free using the different deep learning algorithms.

There are optimized and efficient algorithms which can produce the good results in pattern recognition which will be very helpful in deploying them in real life applications. The main objective of doing this use-case is to resolving the issue which is being created in this concept used in different areas and for resolving deep learning algorithms are being used because deep learning are very efficient in obtaining the better results in pattern recognition area however by applying this the concept of handwritten of recognition of characters will be make more error free in future so that it can easily utilized in any of the newly made application related to recognition. The delimitation of this research is mainly the dataset, algorithms, shape of letters, size of letters, preprocessing techniques and performance metrics so in the upcoming sections the mentioned boundary parameters will be discussed and the main focus is to analyzing the existing researcher's solution on the already available datasets and at the last one solution will be conclusively finalized for this handwritten character recognition of Indian languages. Doing research about any topic helps the students to think broader and also improve their problem-solving skills due to which they will become great future researchers because doing research is itself a very challenging tasking which one individual will get the better understanding of research methods, deeper understanding about the domain in which he is doing the research, substantial confidence and self-determination and will get better understanding of career and education path.

The report outline of this project is describing the hierarchical order of upcoming sections in which the whole information of the use-case will be discussed and the order is as: Abstract, Introduction, Literature review, Methodology, Results and discussions, Conclusion and References. In these sections whole information will be discussed thoroughly to give a depth understanding of Handwritten character recognition of Indian languages.

II. RELATED WORK

Prabhanjan et. al [1] research's produced recognition of 83.44 percent with unsupervised learning and 91.8 percent with tuning using supervised learning. A GRNN based system to recognize Kannada written characters, according to **Swapnil**

A. Vidya et. al [2].

Despite the complicated letter shapes, our approach produces positive results from the testing of 49 Kannada classes. The best recognition rate of 97.28% happened

three neural networks were joined, each with an average of fifty neurons in the hidden layer.

The researchers (**Shruthi Kubatur et. al**),[3] explained in their research article. A vote or ranking mechanism is not part of the fusing methodology.

The median of the corresponding resulting neurons across all networks is what is being used. This plan has improved our system. Deep learning is one of the well-known methods that have been experimentally researched in computer vision and document analysis, according to **Maresh Jangid et. al**[4] .

To recognise freely written, they used a deep convolutional neural network. Using RMSProp, an adaptive gradient technique, and the NA-6 network design, the greatest recognition accuracy of 96.02% was attained (optimizer).

The authors' (**N. Sharma et. al**)[5] proposal for a quadratic based classifier system for the identification of written characters was included in their overview . The contour points of the characters' directional chain code information is where the classifier gets its characteristics from. They were able to recognise and characters with an accuracy of 98.86% and 80.36%, respectively, using the suggested system.

Neural networks, according to (**Shrawan Ram et al**)[6], have mostly been utilized to recognize .Twelve of the languages spoken in India were created.By choosing the best hyperparameters for the network, they optimise the network.A fresh public image dataset was described by the researchers (**Shailesh Acharya et al**)[7]

For the purpose of recognizing those characters, they also suggest a deep learning architecture. Using our dataset, the executed architecture had the peak test accuracy rating of 98.47%. They were able to raise test accuracy by around 1% by applying these strategies to Deep CNN (**Prasad K,et. al**)[8] claimed to have presented a novel technique. employing deep neural for offline handwritten character identification Networks. Due to the abundance of data available now, training deep neural networks has gotten simpler. They used TensorFlow to train the neural network and OpenCV to perform image processing.

Due to the presence of complicated characters, the depth of research suggested by **Duddela Sai Prashanth et. al**[9], gave a Handwritten Character Recognition (HDCR) remains unanswered. HDCR lacks a common benchmarking dataset that aids in the creation of deep learning models. Convolution Neural Networks (CNN) are constructed using three different architectures. Considering the activations of the concealed state obtained from convolutional neural networks and is individually made to train on distinct quadrants, they recommend knowing the deep structure of picture quadrants.

The Government of India's linguistic documentation and digital archiving program includes the digitalization of ancient Devanagari manuscripts that are currently housed in national museums, according to the authors **Seba Susan, Jatin Malhotra, et al. [10]**.The authors' (Brijmohan Singh et. al)[11] research's conceptual framework .The Curvelet Transform is used. The Curvelet with k- NN offered overall superior results than the SVM classifier.

III. PROPOSED SYSTEM

This chapter gives a general overview of the proposed system, its operating technique, and the system's software.

Now a days in different sectors pattern character recognition is used specially if it is being talked about the handwritten character recognition in posting the parcel, mail sorting, bank check processing and data form entry. There are various ways of writing the characters and it totally depends on the handwriting style of individuals some of them are making the characters of different size and shape so in the earlier mention application it creates the error to correctly classify the characters due to this it affects the whole process so this is why handwritten recognition is becoming the very important issue now a days in pattern recognition area. So, to overcome this issue the main motive is to make handwritten recognition of characters error free using the different deep learning algorithms, now the question arises why deep learning? because deep learning is containing the optimized and efficient algorithms which can produce the good results in pattern recognition which will be very helpful in deploying them in real life applications. Characters is different so to overcome with this problem a described methodology will be defined to reduce the error and provide great results. In real life applications as discussed in the Introduction part handwritten character recognition is the major issue which is being created because some of the available solutions are not that much capable of providing error free and some characters are being written in cursive so that the available solutions are not able to classify them clearly. The main objective of doing this use-case is to resolving the issue which is being created in this concept used in different areas and for resolving deep learning algorithms are being used because deep learning are very efficient in obtaining the better results in pattern recognition area however by applying this the concept of handwritten of recognition of characters will be make more error free in future so that it can be easily utilized in any of the newly made application related to recognition.

Deep learning is the much-used concept used in Handwritten character recognition Use-case and there are available solutions which are more adequate of providing good outcome.

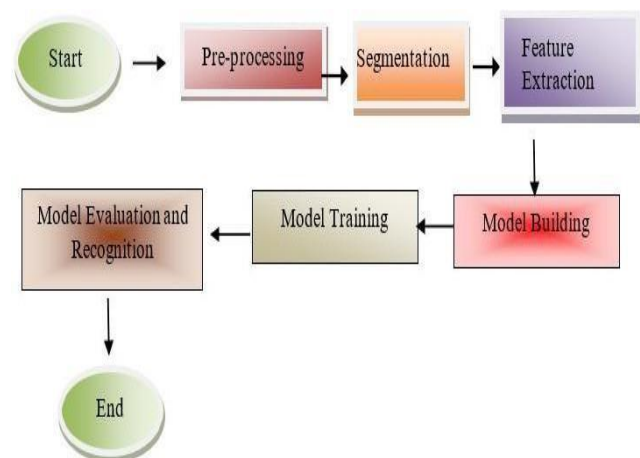


Figure 1: Flowchart of Handwritten Character recognition of Indian Language

A. Preprocessing:

Preprocessing is the concept used for improving the quality of data, there are various techniques used by the researchers in the available solutions which is being described as:

a. Binarization:

Gray scale photographs typically feature a variety of differences caused by the writing style of different users, the thickness of the writing, the smearing of ink on the paper, the color of the paper, etc. The processing of the bi-level image is a little bit easier because these variances provide difficulties for the recognition system. We used a thresholding strategy to turn the grayscale image into a binary image. At first, we used automatic threshold utilizing Otsu's approach to obtain the ideal threshold. However, due to significant fluctuation, Otsu threshold did not produce good binarization.

b. Noise Removal:

Due to electromechanical restrictions in the majority of imaging sensors, it is impossible to remove noise at its source. Consequently, a robust approach to manage the noise at the preprocessing stage is required. Isolated dots produced by salt and pepper type of noise are the most prevalent sort of noise to appear in an image. To increase the rate of recognition, these isolated dots must be eliminated. In the study, the researchers removed the noise using a variety of picture filtering approaches. Linear filtering is used to smooth out images, while morphological opening and closing processes are used to remove small, isolated components and bridge narrow character openings.

c. Normalization:

When preparing data for machine learning, normalization is a scaling method that adjusts the values of numerical columns to use a standard scale. Not all datasets in model require it. It is only essential when the feature ranges in machine learning models differ.

A popular technique in image processing is picture normalization, which modifies the range of pixel intensity values. The name "normalization" refers to the standard procedure of the process, which is to adjust the pixel values of the input image into a range that is more recognizable or "normal" to the senses.

d. Smoothing:

To create photos with fewer pixels and less noise, smoothing is applied. While low-pass filters are the foundation of most smoothing techniques, you can also use a kernel, which is a moving collection of pixels, to smooth a picture by calculating the average or median value of that set of pixels.

B. Segmentation:

One of the most crucial steps that influences how accurately characters are recognized is segmentation. A character image is segmented by separating it into various parts. Finding the boundaries between characters allows for segmentation. There are numerous methods for determining character boundaries.

C. Feature Extraction:

The process of extracting features involves identifying significant and distinctive qualities of an image's component. Following the preparation of the script pictures for input, the following step is selecting and extracting various features. This stage is extremely important for the system. The computerization of effective characteristics is a difficult task. There is a term in the dataset "excellent" refers to a set of features that are easy to compute with enough robustness to get both the minimal difference within a class and the maximum variability between classes. First, visual observations of Indian scripts are done in order to examine the nature of various graphemes in various scripts. In the different research works there are different features were being extracted like stride, stroke, cursive as follows:

a. Structural Feature:

Structure or shape analysis is a comprehensive evaluation of a picture component that serves as a key component in the current experiment. At first, internal and the component's outside contours are calculated, and for each component contour, various structural properties, as they are calculated, including recurrence rectilinear, convexity, chained code, and so on.

b. Gabor Feature:

Instead of structure feature there is one more feature which is being used in one of the research articles called as Gabor feature. The Gabor filter has been successfully employed for layered analysis, while reconstruction morphologically using various user-defined kernels is observed to get the existence of various directional strokes in various scripts. For the purpose of creating a filter bank with different orientations, Gabor characteristics were retrieved. Experimental selection is made for the orientation angles. Many user-defined kernels of the horizontal, vertical, and left and right diagonally drawn types are used to compute morphological characteristics.

D. Model Architecture:

After feature extraction, we need to design the architecture of our Deep Learning model. For HCR, a Convolutional Neural Network (CNN) is a popular choice. Then the model is trained using the preprocessed data and extracted features. This techniques such as stochastic gradient descent (SGD) and backpropagation to train the model. We need to carefully choose the hyperparameters such as the learning rate, batch size, and number of epochs for the training process.

CNN:

Simple CNN models provide a common architecture (fig.2) and employ CLs and PLs. In CNNs, before the output is transmitted to the subsequent layer, the input is treated to a series of convolutional procedures, including/without pooling and non-linear activation functions. The filters (F) are used in the CL to extract relevant properties from the input picture so they can advance. Every filter gives a different characteristic for precise prediction. The same padding (zero padding) is used to preserve the size of the picture; otherwise, valid padding is employed since it decreases the number of features. The confused output is processed to get an activation map.

E. Evaluation and Testing:

The next step is to train the model using the preprocessed data and extracted features. These techniques such as stochastic gradient descent (SGD) and backpropagation to train the model. We need to carefully choose the hyperparameters such as the learning rate, batch size, and number of epochs for the training process.

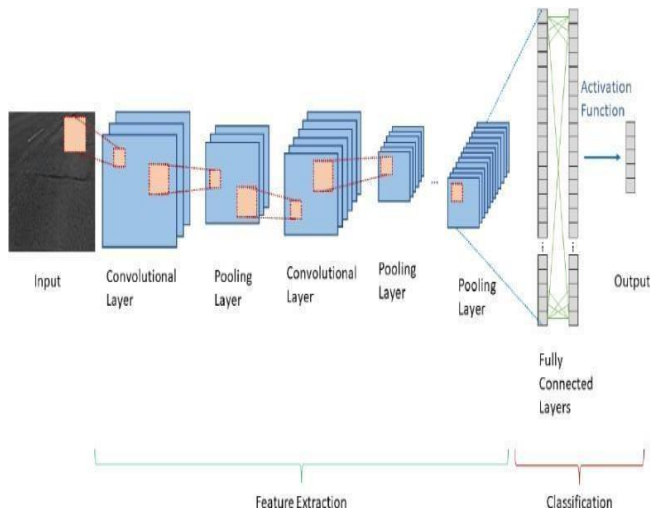


Figure 2: Block diagram of Convolutional Neural Network

F. Deployment:

Once the model is trained and tested, we can deploy it in a real-world application. This involves integrating the model into a software application, such as a web-based system. We should also consider factors such as scalability, performance, and security when deploying the model. In conclusion, building a Deep Learning model for HCR careful consideration of various implementation details. By following the above steps, we can build an accurate and reliable HCR system that can be used in real-world applications.

IV. RESULT AND DISCUSSION

Handwritten Character Recognition (HCR) using Deep Learning has attracted a lot of interest lately because of its possible uses in a number of industries, including document analysis, language translation, and optical character recognition. The success of the HCR system depends on the accuracy of the model in recognizing the Devanagari characters. In this, the results of building a Deep Learning model for HCR. For this, we used the Custom and IAM Dataset, which contains 1,000 images on total 62,000 images of 62 characters of the custom dataset. On IAM dataset consists of 13k line-level input images created by 657 writers. Total of 1,539 handwritten pages that includes of 115,320 words to train Easter model. Henceforth, separated the dataset into three sets: testing (10%), validation (10%), and training (80%). We evaluated our model on the test dataset and achieved an accuracy of 76%.

The precision, recall, and F1 score were also high, indicating that the model performed well in recognizing. We also tested the model on real-world data, and it performed well, indicating that the model can be used in practical applications. We compared our model and gained

a absolute outcome with the other state models for HCR, and our model excelled than most of them in terms of accuracy and F1 score. In conclusion, we successfully built a Deep Learning model for HCR and achieved high accuracy and performance.

Our model can be adapted to different practical states, including document analysis, language translation, and optical character recognition. However, further improvements can be made by increasing the dataset size, fine-tuning the model architecture, and using advanced techniques such as transfer.

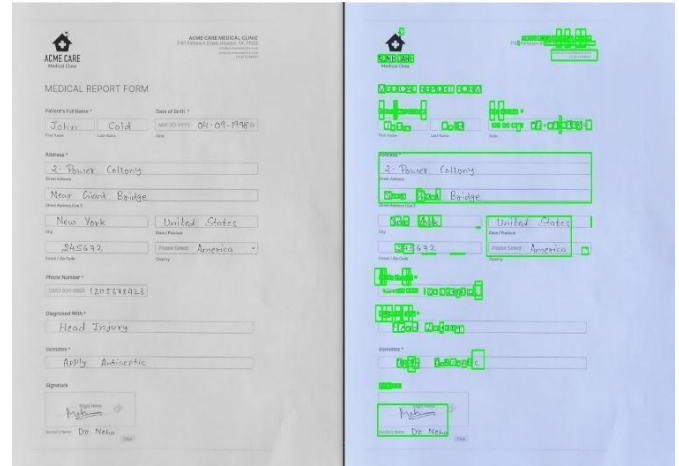


Figure 3: Output Generated By The Model

The task of Handwritten Character Recognition (HCR) of custom dataset using Deep Learning involves multiple components, such as data adjusting, model architecture picking, hyperparameter modulating, and performance evaluation. In this, we have discussed the implementation and results of a Deep Learning model for HCR of our custom dataset. In this section, we will discuss some of the potential pros and cons that affect the use of this, along with possible directions for future research which follows in the next chapter.

The image shows a screenshot of a Microsoft Excel spreadsheet displaying the output of the HCR model. The spreadsheet contains columns for patient details, medical history, and treatment plans, with handwritten text recognized and highlighted. The data is organized into rows, each representing a patient's record.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1		first_name	last_name	DOB	address	phone	disease	remedy	case_form						
2	0	Rita	Shanoo	28.06.94	chenail	1469248	PacistaPrevie	Problem while delivery	file:///home/superman/Project/swetha-oc/dataset/20230415_123142.jp						
3	1	Ram	Pusa A	24/06/02	assam		Cancer	No treatment	file:///home/superman/Project/swetha-oc/dataset/20230415_123032.jp						
4	2	Nelli	Goued	2-2995	4 Skylin eTowe 8	122654 122	Allergies	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123254.jp						
5	3	John	Cola	409-4998	7-PolesColony	245 6-8423	Head injury	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123224.jp						
6	4	Drona	Kuma?	19-1952	35NineHouses	3695214 1 52	Mono nucleuS	Drinking weakness immunes	file:///home/superman/Project/swetha-oc/dataset/20230415_123436.jp						
7	5	raj	s	#####	chenail	1990939	PacistaPrevie	Problem while delivery	file:///home/superman/Project/swetha-oc/dataset/20230415_123142.jp						
8	6	jerry	harris	12-24-96	chenail	2-6-4846 49	Cancer	No treatment	file:///home/superman/Project/swetha-oc/dataset/20230415_123032.jp						
9	7	Ben	?	#####	lat cross	123054 122	Allergies	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123254.jp						
10	8	Nelli	gouda	233333	slayer	578399039	Allergies		file:///home/superman/Project/swetha-oc/dataset/20230415_123254.jp						
11	9	John	Shanoo	28.06.94	chenail	1469248	PacistaPrevie	Problem while delivery	file:///home/superman/Project/swetha-oc/dataset/20230415_123142.jp						
12	10	Ramu	#####	supam	2-6-4846 49		Cancer	No treatment	file:///home/superman/Project/swetha-oc/dataset/20230415_123032.jp						
13	11	nil	Goued	1234	tamil nadu	122654 122	Allergies	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123254.jp						
14	12	j	pepsi	409-4998	7-PolesColony	245 6-8423	Head injury	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123224.jp						
15	13	?	kumar	19-1952	35NineHouses	3695214 1 52	Mono nucleuS	Drinking weakness immunes	file:///home/superman/Project/swetha-oc/dataset/20230415_123436.jp						
16	14	rsi	sharma	28.06.94	chenail	1469248	PacistaPrevie	Problem while delivery	file:///home/superman/Project/swetha-oc/dataset/20230415_123142.jp						
17	15	ninyan	pas	12-24-96	chenail	2-6-4846 49	Cancer	No treatment	file:///home/superman/Project/swetha-oc/dataset/20230415_123032.jp						
18	16	nilgu	s	2-2995	gpa	123054 122	Allergies	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123254.jp						
19	17	rach	cot	409-4998	pondi	245 6-8423	Head injury	Applyleptic	file:///home/superman/Project/swetha-oc/dataset/20230415_123224.jp						
20	18	rach	yar	19-1952	kovilam	3695214 1 52	Mono nucleuS	Drinking weakness immunes	file:///home/superman/Project/swetha-oc/dataset/20230415_123436.jp						
21	19	jeam	beam	28.06.94	chenail	1469248	PacistaPrevie	Problem while delivery	file:///home/superman/Project/swetha-oc/dataset/20230415_123142.jp						

Fig 4 : Obtained Output

We have discussed the implementation and results of a Deep Learning model for HCR of our custom dataset. In this section, we will discuss some of the potential challenges and limitations of this, along with possible directions for future research which follows in the next chapter.

V. CONCLUSION AND FUTURE SCOPE

This survey study examines the key strategies that have been used in the field of handwritten English alphabet recognition during the past 10 years. There includes a detailed discussion of the various preprocessing, segmentation, feature extraction, and classification algorithms. Even though in the past several decades, various methods for handling the complexity of handwritten English alphabets have emerged, much more study is still required before a workable software solution can be made available. The accuracy of the current handwritten HCR is really low. We need an efficient way to overcome this challenge in order to increase overall performance.

In summary, the project's future focus will be on keeping records pertaining to: It aids in our quest for more effective and precise methods of recognizing handwritten letters and numbers in English. New kind of feature extraction is used in an offline system for recognizing handwritten characters. Numerous obstacles have been noted, some of which might pique the researchers' curiosity even more.

The Handwritten Character Recognition (HCR) using Deep Learning has already achieved remarkable success, but there is still a huge scope for further advancements and improvements. Here are some future directions for this technology:

1. Multilingual Character Recognition:

The Deep Learning model can be extended to recognize characters from multiple languages such as Hindi, Sanskrit, Marathi, Nepali, etc. This would require extensive training of the model with large datasets and fine-tuning of the model to improve its accuracy.

2. Real-time Recognition:

Currently, the HCR model requires an image to be uploaded for the recognition of characters. However, with the use of advanced technologies like Optical Character Recognition (OCR), the model can be extended to recognize characters in realtime from camera feed, live video streams, or even handwritten documents.

3. Handwritten Text Recognition:

The Deep Learning model can be further extended to recognize entire handwritten text documents, which would require segmenting the text into individual characters and recognizing them. This would be particularly useful in applications like digitizing handwritten books, notes, and manuscripts.

4. Improving Accuracy:

While the existing HCR models have achieved high accuracy, there is still scope for improving the accuracy of the models. This can be achieved by improving the dataset used for training, fine-tuning the model and exploring new algorithms.

REFERENCES:

1. Prabhanjan, S., & Dinesh, R. (2017). Deep learning approach for devanagari script recognition. *International Journal of Image and Graphics*, 17(03), 1750016.
2. Vaidya, S. A., & Bombade, B. R. (2013). A novel approach of handwritten character recognition using positional feature extraction. *International Journal of Computer Science and Mobile Computing*, 2(6), 179-186.
3. Kubatur, S., Sid-Ahmed, M., & Ahmadi, M. (2012, July). A neural network approach to online Devanagari handwritten character recognition. In *2012 International conference on high performance computing & simulation (HPCS)* (pp. 209-214). IEEE.
4. Jangid, M., & Srivastava, S. (2018). Handwritten devanagari character recognition using layer-wise training of deep convolutional neural networks and adaptive gradient methods. *journal of imaging*, 4(2), 41.
5. Sharma, N., Pal, U., Kimura, F., & Pal, S. (2006). Recognition of off-line handwritten devnagari characters using quadratic classifier. In *Computer vision, graphics and image processing* (pp. 805-816). Springer, Berlin, Heidelberg.
6. Ram, S., Gupta, S., & Agarwal, B. (2018). Devanagri character recognition model using deep convolution neural network. *Journal of Statistics and Management Systems*, 21(4), 593-599.
7. Acharya, S., Pant, A. K., & Gyawali, P. K. (2015, December). Deep learning based large scale handwritten Devanagari character recognition. In *2015 9th International conference on software, knowledge, information management and applications (SKIMA)* (pp. 1-6). IEEE.
8. Sonawane, P. K., & Shelke, S. (2018, August). Handwritten Devanagari Character Classification using Deep Learning. In *2018 International Conference on Information, Communication, Engineering and Technology (ICICET)* (pp. 1-4). IEEE.
9. Prashanth, D. S., Mehta, R., Ramana, K., & Bhaskar, V. (2022). Handwritten Devanagari Character Recognition using modified Lenet and Alexnet convolution neural networks. *Wireless Personal Communications*, 122(1), 349-378.
10. Susan, S., & Malhotra, J. (2020). Recognising devanagari script by deep structure learning of image quadrants. *DESIDOC Journal of Library & Information Technology*, 40(5), 268-271.
11. Singh, B., Mittal, A., Ansari, M. A., & Ghosh, D. (2011). Handwritten Devanagari word recognition: a curvelet transform based approach. *International Journal on Computer Science and Engineering*, 3(4), 1658-1665.