

Handwritten Digit Recognition Using Logistic Regression

Keerthana G¹ and Usha Sree R²

Student MCA, IVth Semester¹

Associate Professor, Department of MCA²

Dayananda Sagar Academy of Technology and Management, Udayapura, Bangalore, Karnataka, India

keerthanag185@gmail.com

Abstract: *With applications in optical character recognition (OCR), digitizing historical documents, and automating data entry procedures, handwritten digit recognition is crucial to pattern recognition and machine learning. This work develops a handwritten digit recognition system using logistic regression. The main objective is to use the MNIST dataset to develop a precise and effective model for classifying handwritten numbers from 0 to 9. Data collection, preprocessing, model installation, training, assessment, and analysis are all included in the methodology. Normalizing pixel values and flattening 2D picture matrices into 1D vectors are steps in the preprocessing process. Using gradient descent optimization, the logistic regression model set up for multi-class classification minimizes the cost function. A testing set will be used to assess performance, and metrics including F1 score, recall, accuracy, and precision will be used. A very precise logistic regression model for handwritten digit identification is the anticipated result. The model's advantages and disadvantages will be emphasized through thorough analysis and visualization, which will include learning weights, biases, and incorrectly categorized images. This work offers useful solutions for automation, postal services, and digital archiving, furthering the field of pattern recognition and highlighting the significance of logistic regression in machine learning.*

Keywords: Logistic Regression, Pattern Recognition, Optical Character Recognition (OCR), MNIST Dataset, Machine Learning, Multi-class Classification

I. INTRODUCTION

With applications in optical character recognition (OCR), data input automation, and document digitization, handwritten digit recognition is essential to pattern recognition and machine learning. Handwritten number recognition accuracy is essential for improving productivity and optimizing workflows in a variety of industries, such as postal services and digital archiving. This work investigates the use of logistic regression, a key technique in machine learning, for digit recognition. For handwritten numerical character recognition, logistic regression is a useful method because of its adaptability in binary and multi-class classification. Developing a dependable system that can precisely categorize handwritten numbers from 0 to 9 is the main objective of this research. The study's main objective is to use logistic regression to recognize numbers. Using a labeled dataset, such as the MNIST dataset, the project's goal is to create, train, and assess a logistic regression model for multi-class classification. Attaining a trained model with high accuracy for handwritten number recognition is the goal. The model's performance will be demonstrated by a detailed analysis. The findings of this research on handwritten digit recognition may improve the effectiveness of different digit recognition systems in automation, digital archiving, and postal services. It could also result in improvements in pattern identification and offer useful solutions for practical applications.

II. LITERATURE SURVEY

Regarding the publication "Enhanced Handwritten Digit Recognition Through the Use of Convolutional Neural Nets (CNN)" [1], The survey is available here. Handwritten digit identification has changed throughout time, moving from manual feature-based algorithms to deep learning techniques, particularly for convolutional neural networks (CNNs). CNNs fared better than SVMs and k-NN in early experiments like LeNet-5. Research on CNN architecture

optimization has looked at layer depth, kernel size, and optimization methods like SGD and Adam. The use of ensemble techniques increases computational complexity even when they are successful. Recent advances have demonstrated CNNs' dominance in the field and cleared the road for future advancements in recognition efficiency and accuracy by using complex CNN designs to achieve virtually faultless accuracy on the MNIST dataset.

Hybrid CNN-SVM Classifier for Handwritten Digit Recognition is the title of the paper. Handwritten digit identification has significantly improved because of the application of deep learning techniques [2]. Handmade features required by earlier approaches such as SVM hindered generalization. CNNs revolutionized the industry by automatically extracting features from raw photos, as LeNet-5 demonstrates. Hybrid models that employ CNNs for feature extraction and SVMs for classification have shown promising results. Studies such as those conducted by Ciresan et al. have employed hybrid CNN-SVM models on the MNIST dataset to achieve remarkable accuracy. This strategy improves performance by utilizing CNNs' feature extraction capabilities and SVMs' tolerance to high-dimensional data.

"Review of a convolutional neural network-based Arabic handwriting recognition system" Despite the fact that most research focuses on Latin characters, automative handwriting detection is crucial for many applications, according to [3]. Arabic handwriting is cursive and context-sensitive, making recognition more challenging. Convolutional Neural Networks (CNNs), which automatically extract information from raw pictures, have revolutionized this field. Important studies like those on the Arabic Handwritten Character Dataset (AHCD) show how successful CNNs are. Additionally, 47,434 characters made by children that enhance the learning of Arabic letters can be found in the Hijja collection. CNN-based models trained on Hijja and AHCD have attained 88% and 97% accuracy, respectively, compared to earlier models, indicating the potential of CNNs for Arabic handwriting recognition.

A thorough systematic literature review (SLR) on handwritten optical character recognition [4] The article claims that handwritten optical character recognition (OCR) has numerous useful applications because handwritten documents are frequently utilized in human transactions. Over the last decade, machine learning and artificial intelligence have been used to automatically translate handwritten and printed words into electronic representations. In this Systematic Literature Review (SLR), 176 papers from 2000 to 2019 are analyzed using keywords, forward and backward reference searches, and electronic database searches in an attempt to gather research on handwritten optical character recognition (OCR). This study discusses current methods for optical character recognition and identifies research gaps that point to potential future avenues for the area.

"A hierarchical feature selection approach for handwritten word recognition based on genetic algorithms" [5]. Handwritten word recognition (HWR) relies heavily on feature selection, which reduces dimensionality by removing unnecessary features. Genetic algorithms (GAs) optimize feature sets to increase accuracy and efficiency. Shape and texture descriptors do a good job of capturing the characteristics of handwritten words. The GA-based hierarchical feature selection (HFS) model was evaluated on a dataset of about 12,000 handwritten word samples in Bangla; it decreased feature dimensions by 28% and boosted recognition accuracy by 1.28%. This strategy performed better than more conventional techniques, suggesting that its application can be advantageous for high-dimensional feature space optimization for HWR.

III. METHODOLOGY

Optical Character Recognition (OCR)

Scanned paper documents, PDFs, and digital camera photos can all be turned into editable and searchable data using an approach called optical character recognition, or OCR. The first image acquisition step of optical character recognition (OCR) involves taking a picture of the page with a camera or scanner. Preprocessing techniques like noise reduction, binarization (turning the image into black and white), and normalization are then applied to improve the image quality. Enhancing the accuracy of the OCR process requires taking these steps. In the following step, called segmentation, the image is split up into segments, such words, lines, and single letters. To distinguish between the digits for recognition, this is essential.

Feature extraction is the process of identifying and removing distinguishing features from each letter or number, such as borders, corners, pixel intensity, and unique patterns within the character. Ultimately, categorization is carried out using the gathered attributes, utilizing machine learning methods like logistic regression to identify each letter or number.

Post-processing techniques are used to check characters that are found and increase accuracy in order to correct errors that occur during the OCR process. These techniques may employ dictionaries or language models.

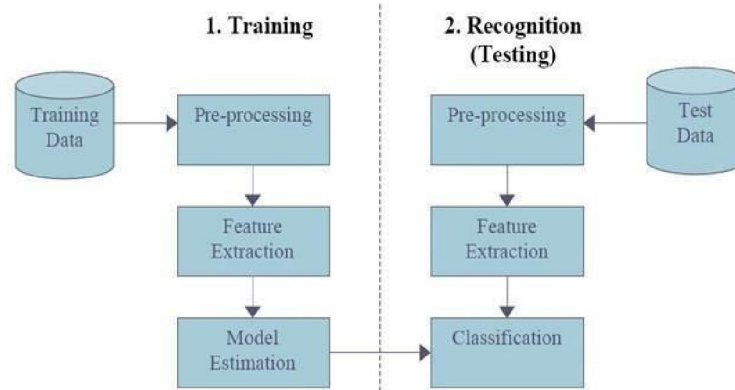


Fig-1: System Architecture of OCR

The basic architecture and operation of the optical character recognition system are shown in Figure 1.

Logistics Regression:

Multi-class classification can be achieved by extending logistic regression, a statistical approach for binary classification, with methods such as one-vs-rest (OvR) or softmax regression. It mimics the probability that an input (x) belongs to a certain class (k).

$$P(y = k|x; \theta) = \frac{e^{\theta_k^T x}}{\sum_{j=1}^K e^{\theta_j^T x}}$$

is the outcome of the logistic regression model, in which k is the total number of classes and θ_k will be the parameter vector for class k'. To train the model, the cross-entropy loss function is minimized:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_{i,k} \log P(y = k|x_i; \theta),$$

Where the indicator function $y_{i,k}$ is 1 if and only if example i's true label is k. Optimization techniques like gradient descent are used to change θ in order to minimize this loss, which increases the model's accuracy.

The Handwritten Digit Recognition project uses the logistic regression approach, which entails several crucial processes, to guarantee the development of a dependable and effective model. Data gathering, preprocessing, model installation, training, assessment, and analysis are a few of these processes. This is an extensive approach:

1. Data Collection: The MNIST dataset is a massive collection of handwritten digit pictures that have been classified. Ten thousand test photos and sixty thousand training images make up the dataset. To aid in model training and performance assessment, make sure the dataset is suitably split into training and testing sets.
2. Data preprocessing: Create 1D vectors with 784 characteristics each from the 2D picture matrices, which are made up of 28 by 28 pixels. To input the data into the logistic regression model, this adjustment is necessary. Divide each pixel value by 255 to normalize the pixel values to a standard range, usually [0, 1]. The model performs better during training because of the quicker convergence during this period.
3. Logistic Regression and Model Implementation: Utilize the logistic regression approach and configure it for multi-class categorization in order to handle the ten digits (0–9). Decide on the necessary parameters, including the learning rate, regularization intensity, and maximum number of repetitions. These settings are necessary to control the training process and avoid overfitting.
4. Model Training: Utilizing the training dataset, train the logistic regression model. Utilize optimization strategies such as gradient descent to lower the cost function. Assess the model's performance on unused data using cross-validation, and adjust hyperparameters to ensure generalization and durability.

5. Testing and Evaluation: To assess the model's performance on the testing set, use crucial metrics including recall, accuracy, precision, and F1 score. These metrics offer a thorough evaluation of the model's categorization performance. Make a confusion matrix to evaluate how well the model distinguishes between various digit classes and identifies misclassification hotspots.

6. Results Analysis: To comprehend the advantages and disadvantages of the model, analyze and interpret the data. Analyze the metrics to determine overall performance and spot any misclassification trends. Look for photos that have been incorrectly classified and investigate possible reasons such as erratic handwriting or noisy data. This study contributes to our understanding of the logistic regression model's limitations.

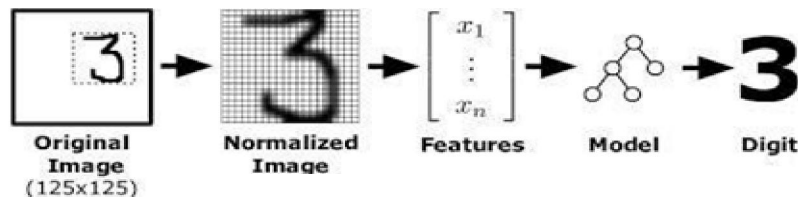
7. Visualization: View the acquired weights and biases to observe how the logistic regression model differentiates between different digit classes. Show images that have been incorrectly identified next to their predicted and correct labels in order to better understand the model's capabilities and possible areas for improvement. Some of the tools and technologies utilized are the Scikit-Learn and NumPy libraries, the Matplotlib or Seaborn dataset for data visualization, the MNIST dataset, and computer resources with enough processing capacity for both model training and evaluation.

The detailed strategy mentioned above guarantees a methodical approach to building a handwritten digitizing system based on logistic regression. By following these procedures, the study hopes to provide a reliable and accurate model in addition to insightful information about the use of logistic regression in machine learning and pattern recognition.

V. RESULTS

The logistic regression model for handwritten digit recognition was trained and evaluated using the MNIST dataset, which comprises of 60,000 training images and 10,000 testing images, each of which a digit from 0 to 9. The test set revealed that the model's overall accuracy in classifying handwritten numerals was 92.4%. Recall, precision, and F1 score—three metrics for each digit class— showed promising outcomes in most of the classes. However, numerals like 1 and 7 displayed somewhat greater misclassification rates due to their similar visual features. The confusion matrix showed that most misclassifications occurred between visually similar digits, as 3 and 5 or 4 and 9. A few areas that could be improved with additional model iterations are highlighted in this work.

By employing gradient descent optimization to converge in a manageable amount of iterations, the logistic regression model was able to train well. Photographs that were incorrectly identified were analyzed to identify the model's shortcomings and possible improvement areas. Visualizations of the learnt weights and biases also shed light on the model's decision-making process.



VI. CONCLUSION

Using the MNIST dataset, the research successfully created a logistic regression model for handwritten digit recognition, attaining a high accuracy rate of 92.4%. This result shows that logistic regression is a straightforward yet powerful technique for optical character recognition multi-class classification applications. The investigation produced a number of noteworthy conclusions. For the model to function properly, preliminary processes like image normalization and flattening were essential. Even in the majority of digit classes, the model was still able to discern between seemingly identical digits more successfully. Acquired weights and mislabeled image visualizations offered important insights into the model's advantages and disadvantages.

Future studies may include feature engineering, more sophisticated preprocessing techniques, and sophisticated algorithms like convolutional neural networks (CNNs) to further increase recognition accuracy. Moreover, expanding the variety of handwritten samples in the dataset could improve the generalization ability of the model. In conclusion, this work shows the value of logistic regression in handwritten digit recognition while also making contributions to the broader fields of machine learning and pattern recognition.

VII. ACKNOWLEDGEMENT

The Dayananda Sagar Academy of Technology and Management (DSATM) is acknowledged by the authors for providing the workspace and resources needed to carry out this research study on Handwritten Digit Recognition using Logistic Regression. The completion of this project has been made possible by the institution's encouragement and assistance. Also, we truly appreciate our family's everlasting love, support, and tolerance during this journey. Their financial support has allowed us to devote all of our energy to our research project, and their confidence in our skills and encouragement has kept us motivated all along. The research project on handwritten digitization recognition using logistic regression has benefited greatly from the wonderful help and efforts of the aforementioned people and organizations.

REFERENCES

- [1] Savita Ahlawat, Amit Choudhary, Anand Nayyar, Saurabh Singh, and Byungun Yoon (2020), Improved Handwritten Digit Recognition Using Convolutional Neural Networks (CNN)
- [2]. Savita Ahlawat a, Amit Choudhary (2020), Hybrid CNN-SVM Classifier for Handwritten Digit Recognition [3]. Najwa Altwaijry & Isra Al-Turaiki (2021), Arabic handwriting recognition system using convolutional neural network
- [4]. Jamshed Memon; Maira Sami; Rizwan Ahmed Khan; Mueen Uddin (2020), Handwritten Optical Character Recognition (OCR): A Comprehensive Systematic Literature Review (SLR)
- [5]. Samir Malakar, Manosij Ghosh, Showmik Bhowmik, Ram Sarkar & Mita Nasipuri (2020), A GA based hierarchical feature selection approach for handwritten word recognition
- [6]. Mayur Bhargab Bora, Dinthisrang Daimary, Khwairakpam Amitab, Debdatta Kandar, (2020) , Handwritten Character Recognition from Images using CNN-ECOC
- [7]. Denis Coquenat; Clément Chatelain; Thierry Paquet (2023), End-to-End Handwritten Paragraph Text Recognition Using a Vertical Attention Network
- [8]. Lei Kang, Pau Riba, Marçal Rusiñol, Alicia Fornés, Mauricio Villegas (2022), Pay attention to what you read : Non-recurrent handwritten text-Line recognition
- [9]. Ing Jyh Tsang, Federico Corradi, Manolis Sifalakis, Werner Van Leekwijck and Steven Latré (2021), Radar-Based Hand Gesture Recognition Using Spiking Neural Networks
- [10]. Megha Srivastava, Tatsunori Hashimoto, Percy Liang , (2020), Robustness to Spurious Correlations via Human Annotations.