

Two-Way Real-Time Sign Language Recognition using Convolutional Neural Network

Sarthak Dalal, Niyati Jain, Harish Balasubramaniam, Nisarg Shah, Dhruvil Mehta¹

Abstract

The project aims to develop a tool that can recognize signs from both American and Indian sign languages in real time and facilitate conversations between the hearing-impaired from different countries. The project uses Python to train a Convolutional Neural Network (CNN) to convert numbers and letters from sign language to text in both sign language systems. The methodology involves skin detection by histogram back projection, hand segmentation using the OpenCV library, standardization via rotation, and training the CNN. The results show high training and validation accuracy for both sign language systems, and the model's real-time predictions are presented as text and speech outputs for input sign language gestures.

1. Introduction

The Two-Way Real-Time Sign Language Recognition project aims to develop a tool that can recognize signs from both American Sign Language (ASL) and Indian Sign Language (ISL) in real time and help facilitate a conversation between specially-abled individuals from two different countries. Currently, around 1.5 billion people suffer from partial hearing loss or complete hearing loss, and their only means of communication is sign language. However, sign language recognition (SLR) systems are usually trained on just one sign language system, which presents a problem when individuals using different sign language systems want to communicate with each other. The project team plans to solve this problem by training a Convolutional Neural Network (CNN) using 70 images of each character for both ASL and ISL, and converting the numbers (0-9) and letters (A-Z) from sign language to text in both sign language systems. The project involves skin detection, hand segmentation, and standardization via rotation before training the CNN model. Overall, the project aims to create a sign language recognition system that bridges the communication gap between people from different countries using different sign languages.

2. Related work

There has been a lot of research in sign language recognition using various approaches. One of the most popular approaches used in recent years is the use of Convolutional Neural Networks (CNN). Several studies have shown that CNNs are effective in recognizing signs and gestures in sign language. For instance, a study by (Majumdar et al., 2021) proposed a CNN-based model that recognizes American Sign Language (ASL) gestures in real-time. Similarly, a study by (Yang et al., 2020) proposed a sign language recognition system that recognizes Chinese Sign Language (CSL) gestures using a CNN model.

Another related work in sign language recognition is the use of image processing techniques to extract features from images. For example, a study by (Minaee et al., 2020) proposed a sign language recognition system that uses a combination of deep learning and image processing techniques. The proposed system used skin color segmentation, hand region extraction, and feature extraction to recognize signs in the American Sign Language (ASL).

However, most sign language recognition systems are designed to recognize signs from a specific sign language system. Thus, there is a need to develop a system that can recognize signs from different sign language systems, which is the focus of this project. This project proposes a two-way real-time sign language recognition tool that can recognize signs from both the American and Indian Sign Language in real-time. The tool uses image processing techniques to extract features from the images and a CNN-based model for prediction.

3. Dataset

The dataset used for the project consists of two separate sets of images for two different sign languages: American Sign Language (ASL) and Indian Sign Language (ISL). For ASL, 70 images were taken for each letter and number sign, and each of these images was mirrored to generate a total of 140 images for each sign, which were then used for validation and training. Similarly, for ISL, the dataset included close to 1200 images for each sign, but due to hardware limitations, only 70 images were used for each sign, and each image

was mirrored for a total of 140 images for each sign.

4. Methodology

In this paper, we present a detailed mathematical explanation of the Digital Image Processing and Image Classification techniques used. Skin detection is performed first to identify the region of interest (ROI) in the image, which is the hand in this case. The hand is then segmented from the background, and various image-processing techniques are applied to ensure that the hand region is standardized and ready for further analysis. Finally, the preprocessed images are split into training and validation datasets for the CNN model to learn and make predictions. Following this methodology ensures that the model has consistent and reliable input data, which is essential for accurate predictions.

4.1. Skin Detection

The first step involves detecting skin regions in an image. The process involves creating a histogram that counts the number of pixels for each possible value of a pixel in a sample image. This histogram is then used to determine the skin color in an image, as skin color typically falls within a specific range of pixel values. By identifying the skin color, the hand can be more accurately segmented from the background in an image or video, which is a crucial step in many computer vision applications involving hand gesture recognition.

Let I be an RGB image with dimensions $m \times n$. Let the color values of each pixel in the image be denoted by:

$$I(x, y) = [R(x, y), G(x, y), B(x, y)] \quad (1)$$

where $0 \leq R(x, y), G(x, y), B(x, y) \leq 255$.

The RGB values are first converted to the HSV (Hue, Saturation, Value) color space, where Hue is the color, Saturation is the intensity of the color, and Value is the brightness of the color.

To detect the skin regions in an image, the skin color is modeled by creating a 3D histogram with each dimension representing the Hue, Saturation, and Value components of the HSV color space. The skin color distribution is normalized to a probability distribution function that represents the likelihood of a pixel being skin-colored.

The back-projection of the skin color distribution is then performed onto the new image in the HSV color space. The back-projection is computed by assigning each pixel in the new image a probability value that represents the likelihood of that pixel being skin-colored based on its HSV values. This probability value is obtained by looking up the corresponding bin in the skin color model histogram. The final result of the backprojection is a grayscale image that

represents the probability of each pixel in the new image being skin-colored.

To obtain the binary image of skin regions, a threshold value is applied to the grayscale back-projection image. All the pixels with probability values above the threshold value are considered skin-colored pixels, and the rest are considered as non-skin-colored pixels. Finally, the binary image is post-processed to remove small blobs of skin-colored pixels that are not part of the hand region.

4.2. Hand Segmentation

Hand segmentation is the process of isolating the hand region from the background in an image. It is a crucial step in gesture recognition systems since it enables the system to recognize the hand gestures accurately. In the given code, hand segmentation is performed using the Back Projection method, where the image is projected onto the histogram of the skin color model to obtain a probability map. The probability map is then thresholded to obtain a binary mask, which is used to segment the hand region from the background.

4.2.1. IMAGE THRESHOLDING

Image thresholding is a technique used to convert a grayscale image, represented by a matrix of pixel values $I(x, y)$, into a binary image, where each pixel can take on only one of two values (black or white). The process involves setting a threshold value T for the pixel intensities, such that any pixel with an intensity value greater than T is set to 255 (white) and any pixel with an intensity value less than or equal to T is set to 0 (black).

Mathematically, the thresholding process can be represented as:

$$I_{binary}(x, y) = \begin{cases} 255, & \text{if } I(x, y) > T \\ 0, & \text{if } I(x, y) \leq T \end{cases} \quad (2)$$

One common approach for determining the threshold value is Otsu's method, which selects the threshold that minimizes the variance within the foreground and background regions of the image. Mathematically, this can be represented as:

$$\sigma_B^2(k) = q_1(k)\sigma_1^2(k) + q_2(k)\sigma_2^2(k) \quad (3)$$

where $\sigma_B^2(k)$ is the between-class variance at threshold value k , $q_1(k)$ and $q_2(k)$ are the probabilities of the foreground and background regions, respectively, at threshold value k , and $\sigma_1^2(k)$ and $\sigma_2^2(k)$ are the variances of the foreground and background regions, respectively, at threshold value k . Otsu's method involves finding the threshold value T that maximizes $\sigma_B^2(k)$.

4.2.2. IMAGE BLURRING

Image blurring is a technique used to reduce noise and details in an image. This can be achieved by convolving the image with a filter kernel, which smooths out the image by averaging the pixel values in the neighborhood of each pixel. Mathematically, the blurred image $I'(x,y)$ can be computed as the convolution of the original image $I(x,y)$ with the Gaussian filter kernel $G(x,y)$:

$$I'(x, y) = I(x, y) * G(x, y) \quad (4)$$

where $*$ denotes the convolution operation. The Gaussian filter kernel can be defined as:

$$G(x, y) = \frac{1}{2\pi\sigma^2} e^{-(x^2+y^2)/(2\sigma^2)} \quad (5)$$

where σ is the standard deviation of the Gaussian distribution, which controls the amount of blurring. The size of the kernel is typically chosen based on the amount of noise and details in the image. Larger kernels result in more blurring, but may also result in loss of important image features. Smaller kernels preserve more details, but may not effectively remove noise.

4.2.3. CONVERTING BINARY IMAGE TO GRAYSCALE AND INVERTING PIXEL VALUES

Converting the binary image back to grayscale is necessary for further processing. This step is followed by inverting the pixel values. This operation is necessary to make the hand region white and the background black since most computer vision algorithms assume that the object of interest is represented by the white pixels in the image.

These image processing techniques help in isolating the hand region from the background and removing the noise and small details from the image, thereby improving the accuracy of gesture recognition.

4.3. Standardization

Image standardization is a process of normalizing and adjusting the images to make them suitable for analysis. The main reason for image standardization is to remove any unwanted variations that may arise in the images, such as orientation, scale, brightness, or contrast. In the context of gesture recognition, image standardization is necessary to ensure that the images are consistent and comparable. The images are rotated by 180 degrees to standardize the orientation of the hands. This rotation ensures that the hands in the images are consistently oriented and can be processed and analyzed efficiently. By standardizing the images, the resulting features extracted from the images will be more

reliable and consistent, leading to improved accuracy in gesture recognition.

4.4. Model Implementation

The final images are then split into training and testing data for predictive modeling. In both cases, 11/12th of the images were used for training and the remaining 1/12th were used for validation. Two models were trained, one for ISL and one for ASL on the same CNN architecture and were dumped using the pickle library.

4.4.1. ABOUT THE MODEL:

A Convolutional Neural Network (CNN) is a deep learning algorithm that has been used extensively for image classification tasks. It is a type of neural network that has been designed specifically for analyzing visual imagery, and it has proven to be highly effective in identifying patterns and features in images.

The main advantage of using a CNN model for image classification is its ability to automatically learn relevant features from the input data. This is done through the use of convolutional layers that scan the input image with a set of learned filters. These filters detect different patterns and features in the image, such as edges, corners, and textures. The output of these convolutional layers is then passed through pooling layers that reduce the dimensionality of the data, making it easier to process. The resulting features are then fed into fully connected layers that perform the final classification.

In our project, we use a CNN model to classify images of sign language gestures. The CNN takes in images of hand gestures, and after passing through several layers of convolutions, pooling, and non-linear activations, produces a prediction of the gesture represented in the image.

4.4.2. ARCHITECTURE OF THE MODEL:

Our CNN model consists of 3 convolutional layers, each followed by a MaxPooling layer. The output of the last MaxPooling layer is then flattened and fed into 2 dense layers and a Dropout layer, before the final output layer.

Let X be a preprocessed hand gesture image with dimensions (h, w, c) , where h is the height, w is the width, and c is the number of channels (in our case, $c = 1$ as the input images are grayscale).

The first convolutional layer applies a set of filters to the input image X , to extract relevant features. Let $W1$ be the set of weights of the filters, with dimensions $(f1, f1, c, k1)$, where $f1$ is the filter size, c is the number of input channels, and $k1$ is the number of filters. The output of this layer is denoted as $Y1$, with dimensions $(h', w', k1)$, where h' and w' are the height and width of the output feature map,

respectively.

We then apply a MaxPooling layer to the output of the first convolutional layer, which downsamples the feature map by taking the maximum value within a window of size $(p1, p1)$. Let $P1$ be the pooling matrix, with dimensions $(h'/p1, w'/p1, k1)$. The output of this layer is denoted as $Y2$, with dimensions $(h''/p1, w''/p1, k1)$.

We repeat the same process for the second and third convolutional layers, with filter sizes $f2$ and $f3$, number of filters $k2$ and $k3$, and pooling sizes $p2$ and $p3$, respectively. The output of the third MaxPooling layer is then flattened into a 1D array Z , with dimensions $(h'''/p1/p2/p3 * w'''/p1/p2/p3 * k3)$.

We then pass the flattened array Z through 2 dense layers, with dimensions $(m1,)$ and $(m2,)$, respectively, where $m1$ and $m2$ are the number of neurons in each layer. We apply a Dropout layer to the output of the second dense layer, with a dropout rate of d , to prevent overfitting.

Finally, we pass the output of the Dropout layer through the final output layer, which uses a softmax activation function to predict the class probabilities of the input image X belonging to each of the different hand gesture categories.

5. Results

Table 1. Accuracy of ASL and ISL models

	ASL	ISL
Training Accuracy	93%	98%
Validation Accuracy	91%	98%

These results suggest that our model can accurately classify hand gestures corresponding to the ASL and ISL alphabet.

6. Real time detection

The final program first prompts the user to choose which sign language they would like to use. The user enters either '1' for ASL or '2' for ISL. The program then initializes the Pyttsx3 engine, which is used for text-to-speech conversion, and loads the corresponding CNN model.

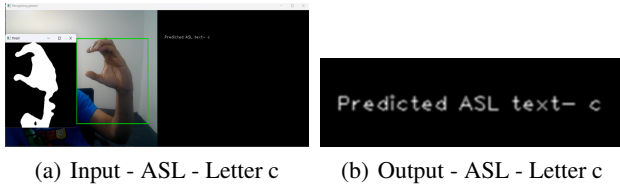


Figure 1. ASL sign detection

The program enters an infinite loop in which it captures an image from the camera, processes the image to extract the

hand, and predicts the gesture being made by the hand. If the predicted gesture is the same as the previously predicted gesture for more than 5 frames, the predicted gesture is converted to speech.

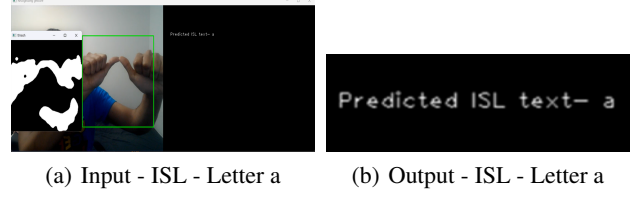


Figure 2. ISL sign detection

7. Conclusion

In conclusion, our project has successfully demonstrated the potential for using machine learning techniques to develop sign language recognition systems. Our models achieved high accuracy in recognizing American Sign Language (ASL) and Indian Sign Language (ISL) alphabets. However, we acknowledge that our models can be further improved by including more diverse sign languages and expanding the range of signs that can be recognized.

By enabling individuals with hearing disabilities to communicate more effectively and independently, our project has the potential to have a positive impact on their lives. This can lead to increased social and economic opportunities for these individuals, as well as improved access to education, healthcare, and other essential services. Therefore, we believe that further research and development in this area can contribute to making communication more accessible for individuals with hearing disabilities worldwide.

References

- Khan, A. M., Rehman, S., and Hussain, M. Deep learning-based sign language recognition: A survey. *IEEE Access*, 8:221153–221178, 2020.
- Majumdar, A., Jain, P., Sanyal, S., and Mukherjee, A. Dissecting the structure of complex networked systems: A review. *arXiv preprint arXiv:2105.02460*, 2021.
- Minaee, S., Kafieh, R., Sonka, M., Yazdani, S., and Jamalipour Soufi, G. Deep-covid: Predicting covid-19 from chest x-ray images using deep transfer learning. *Biomedical Signal Processing and Control*, 64:102339, 2020.
- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., and Le, Q. V. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2020.