

A Recognition Method of CAPTCHA with Adhesion Character

Huo Hua^{1,2} and Chang Guoqin^{1,*}

¹Laboratory of Intelligent Computing & Application Technology for Big Data,
Henan University of Science & Technology, Luoyang, Henan, 471023, China

²School of Software, Henan University of Science & Technology,
Luoyang, Henan, 471023, China

*Corresponding author: 15829715097@sina.cn

Abstract

The emergence of CAPTCHA (Completely Automated Public Turing Test to Tell Computers and Humans Apart) is to better protect the network security ,and the research on recognition of CAPTCHA technology is conducive to expand the design ideas and to improve the loopholes of the original design. In order to improve the recognition rate of CAPTCHA, we proposed a recognition method of CAPTCHA with Adhesion Character; which effectively improves the segmentation quality of adhesive character CAPTCHA in the complex background. First is the based preprocessing of the image and utilize the method of connected area noise reduction to further denoise the image. Then use the projection histogram to cut the adhesion characters and finally neural network training is used to obtain the final recognition results. The method can deal with images with complex background, basically achieve zero noise cutting and can better cut the adhesion characters, the latter part of the training process is relatively simple. The experimental results show that this method is practical and has higher recognition rate.

Keywords: CAPTCHA Recognition; Image Processing; Adhesive Character Segmentation; BP neural network

1. Introduction

With the rapid development of network technology, the function and use of the network is more and more comprehensive and convenient, people also rely more on the network. The network has become one of the important parts of people's work, study and entertainment activities[1]. Meanwhile, people in the process of using the network also produced a lot of privacy information and information as a valuable human resources the security is more and more important. Some criminals through high-performance hardware devices, specific features of malicious programs ,internet worm ,Web design vulnerabilities and other means to crack account passwords, steal user's information and malicious automatic registration, etc. These illegal attacks increased the burden on database and server of the site, resulting in a large number of user's information leakage, affecting the user experience, even lead to the site crash that the site cannot be used normally.

In order to ensure the security of network and strengthen the protection of information, CAPTCHA ,which operates simply and quickly , began to be widely used in the major sites .Compared with other security authentication methods, it contains a small amount of data and also improve the site security and anti-attack capabilities effectively. CAPTCHA is a fully automated Turing test procedure for distinguishing between computers and humans that was first introduced in 2000 by Luis von Ahn *et al* [2].

How to design a secure and efficient CAPTCHA has become a major concern for the

designers of the major web sites. The research of CAPTCHA recognition technology is of great significance to improve the network security performance. In the process we can understand the rule and design principle of the CAPTCHA and helpful to find the defect design, in order to make CAPTCHA become more secure , more mature and more in line with user experience.

In this paper, we proposed a method based on BP neural network to recognize the character adhesion CAPTCHA and the experiment used C++ coded. First of all, through the gray,image corrosion,binarization and other operations on the image preprocessing, followed by the use of projection histogram to cut the characters, then use the BP neural network for training and recognition, and last analysis about the experimental results and existing problems briefly.

The experiment shows that used image corrosion for image denoising can be clear character outline in the recognition method based on the BP neural network, at the same time, choose a reasonable sequence of noise reduction can greatly improve the performance of image denoising, segmentation and projection histogram can better adhesion character, finally training recognition has not less than 92% correct recognition rate.

2. Related Works

As the widely use of the CAPTCHA ,researchers have also designed different types of CAPTCHA for different security[3-4]. Rusu et al. based on the characteristics of the handwriting machine was difficult to recognize. They proposed a text CAPTCHA with handwriting[5-6];Google introduced a new type of CAPTCHA called What's up which was based on image orientation and requires the user to rotate an image to its vertical orientation[7]; At the same time, researchers also designed a sound CAPTCHA based auditory, GAO et al. who were research CAPTCHA recognition on the voice recognition type and jigsaw puzzle type[8-9]. However, these CAPTCHA are generally exist issues like difficult to be identified by users, identification need a long time, taking up storage space and so on, so the text-based authentication code is still the mainstream used by the major sites.

Whether the text-based CAPTCHA could use OCR (optical character recognition) technology to extract the text in the image? OCR could recognize characters in text files very well andearly CAPTCHA can achieve good results by using OCR extraction. However, now there are a lot of noise interference in the CAPTCHA and it is not ideal to use OCR to extract characters directly.WANG et al. used the SVM training to recognize the image and compare it with the OCR directly. The experiments show that SVM training has better recognition effect[10].

For the character segmentation and contour extraction in the image, Shih-Yu et al. proposed a projection-based segmentation method for early CAPTCHA in Yahoo and MSN. The experiment results show that the proposed method had better for character segmentation effect[11]. SHIH et al. proposed the projection and the axis of the segmentation technique to segment the twist CAPTCHA, the experiment results show that the segmentation rate of 75%[12]. JS Cui et al. proposed a new method base-pixel-depth to extract the skeleton of the CAPYCHA,the experiment results show that compared with the traditional skeleton extraction method it reduces the computational complexity and improves the efficiency of image processing[13].

At the same time there were many other ways to recognitionCAPTCHA. GUO et al. through the process on image binarization, de-noising, dilation, splitting characters and then gives out a algorithm based on edit distance which defines the string similarity. Experiments show that the proposed algorithm is simple, fast, robust performance and has a high recognition accuracy rate[14].Zhang L et al. proposed a recognition method based on LSTM model RNNand achieved good results[15]. ZHANG et al. proposed a

new image analysis model named Concept Component Analysis (CCA). This new model based on the approaching idea in Newton's iteration and it was solved by a multi-population genetic algorithm, in their experiments has achieved good recognition results[16]. CHEN et al. proposed the probability pattern framework to recognize the target numbers in the CAPTCHA images and the experiment proved this method achieved an average of 81.05% for more than two thousand CAPTCHA cases[17]. YIN et al. proposed a recognition method based on density-invariant feature transform and random sampling consistency algorithm. Experiments show that this method has good effect on distorted adhesion CAPTCHAs of different difficulty levels[18].

Artificial Neural Network (ANN) as one of the most popular technologies in recent years has also been widely applied to the CAPTCHA recognition technology. Chellapilla and Patrice used machine learning methods to crack several authentication codes from Web sites and the recognition rates ranging from 4.89% to 66.2% [19]. YIN et al. designed a SVM based on the image CAPTCHA recognition and achieved good results[20]. Wang et al. designed a recognition algorithm based on the convolutional neural network used for recognition have some adhesion characters CAPTCHA in the site of Maopu, Xici Hutong and so on, they all had a higher recognition rate[21].

At present, there are a series of problems in recognition CAPTCHA such as the lack of generality, inefficient character segmentation of overlapping, distorting and blocking and so on. For these problems, this paper proposes a character recognition algorithm based on BP neural network which will be described in detail below.

3. Recognition Method for Adhesion Character CAPTCHA

In order to improve the security of the CAPTCHA and reduce its recognition rate, most of the CAPTCHA is not in the form of a single character, but the content contains numbers, letters, characters, complex background or noise, etc. In the process we first need to preprocess the collected CAPTCHA pictures which include part of the denoising, character segmentation and so on. Next, we used neural network for character recognition and verification. Process shows in Figure 1.

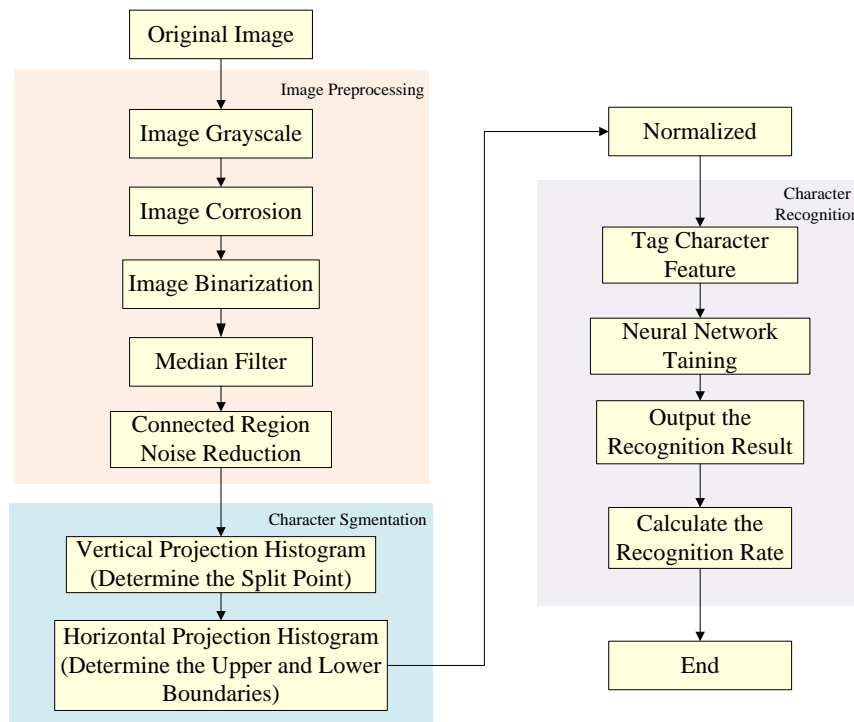


Figure 1. CAPTCHA Identification Process

3.1. Image Preprocessing

Image preprocessing is the most basic step in the CAPTCHA recognition process and the result of the preprocessing will affect the character segmentation and extraction the feature. Therefore, the collected image needs series of operations like image grayscale, denoising and so on , It was purpose make original image become easy to segmentation.

(1) Color Image Gray

Most of the CAPTCHA is a color image with noise that purpose is to increase the difficulty of recognition, confusion characters and background. Color images contain more information and more difficult to deal with^[22], so most of the images are transformed from RGB images to grayscale images without affecting the recognition. The first step of image gray-scale as preprocessing is to lay the foundation for post-image processing. Figure 2-1 shows original image, Figure 2 -2 shows grayscale effect.



Figure 2-1. Original Image



Figure 2-2. Image Grayscale

(2) Image Denoising

Image denoising is the most critical step in the image preprocessing. In the case of keeping the edge of the character clear, we can get rid of most of the noise by corrosion, image binarization and median filtering.

Corrosion: Through the corrosion of the image to reduce the larger noise in the image and remove the smaller noise, but also corrode a part of the adhesion of the edge of the character noise, reducing the noise on the edge of the impact of the characters, and make the character edge more clear. Figure 2-3 shows effect;

Image binarization: In order to ensure that the lighter color in original image to retain, when set the threshold value parameter of the binarization should not be too large, while the image on the contrary to make it black and white characters more likely to highlight the noise contained in image. Figure 2-4 shows effect;

Median filtering: Median filtering is the most commonly used filtering method in image processing. It was found that the noise particles in the image are small similar to the salt and pepper noise it means the median filter can achieve good denoising effect. If at the beginning using the median filter, will lead to background and character blur, character outline is not clear enough and increase the difficulty of late character recognition. Figure 2-5 shows effect;



Figure 2-3. Corrosion



Figure 2-4. Image Binarization



Figure 2-5. Median Filtering

(3) Connected region noise reduction

After the basic de-noising process is finished, the further de-noise processing is needed to ensure that there will be no noisewould effect to the character segmentation results. Markall the connected region in the image that meansall the large and isolated noise will be clearly marked out. Since the required character ensure was a large connected region, use the size of the connected region to set the threshold and the smaller connectedregion was deleted to achieve the purpose of de-noising. By connected region noise reduction, most of the noise in the image has been removed that leaving only the characters to be segmented. After the early denoising use of connected area noise reduction that avoid the noisearound the adhesion characters to calculation into the same connected area, lead to the character border blurred and affecting the subsequent character segmentation. Figure 2 -6shows effect.



Figure 2-6. Connected Region Noise Reduction

3.2 Character Segmentation and Normalization

Character segmentation was the most important part of the CAPTCHA recognition technology. It is necessary to retain the integrity form of the character and the adhesive character will be separated from each other. The single character can be input to the neural network as a test sample, this sample will help extract the character features and facilitate the training and classification behind.

The image after the denoising shows the vertical direction of the grayscale projection histogram, the value of the zero position is the non-stick character character boundary, the relative width of the character is calculated, and the projection histogram can be found between each two characters Obviously the trough uses the character, uses the character width and the trough position to determine the demarcation point of the character, determines the cutoff points and cuts the image according to the demarcation point. After character was separated used the gray-scale projection histogram in horizontal direction to find the character upper and lower boundaries got the single character image. Figure 3-1 and Figure 3-2 shows segmentation effect.

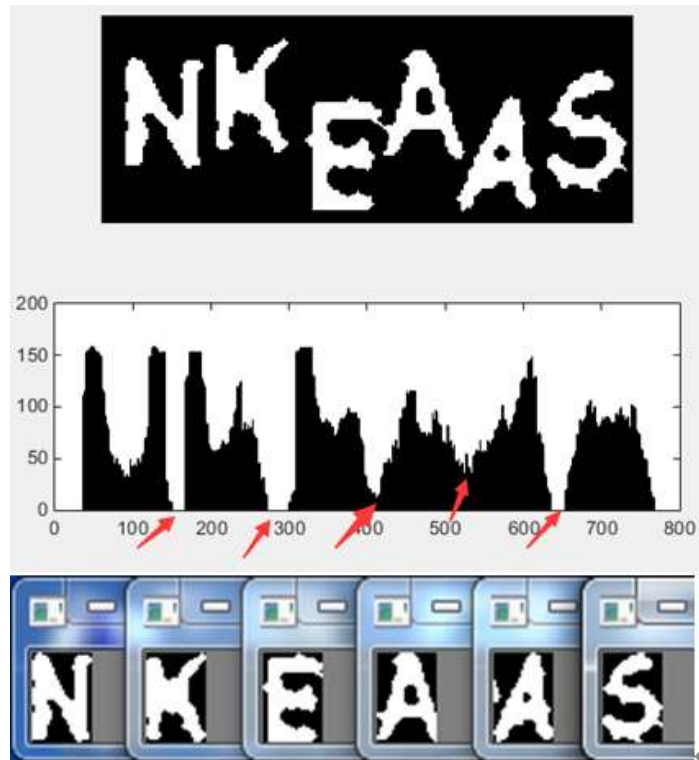


Figure 3-1. Image Segmentation Effect

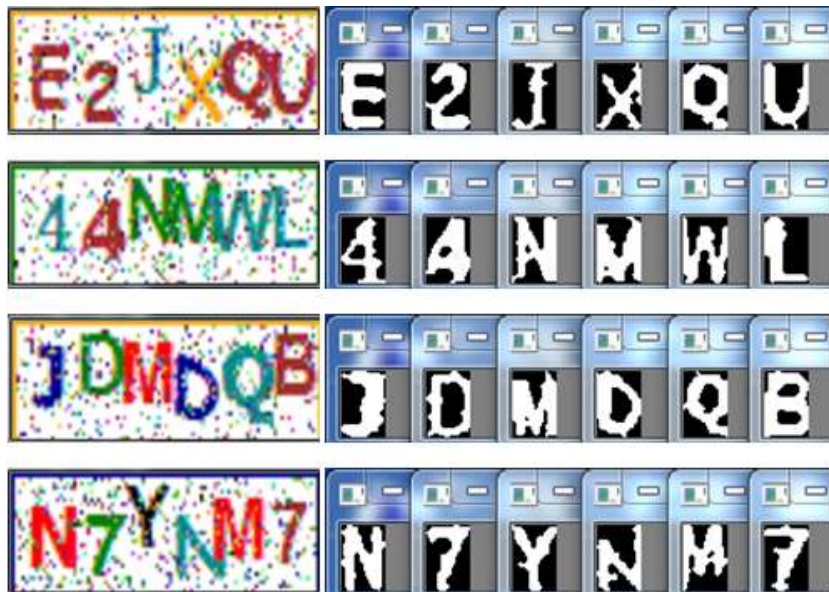


Figure 3-2. Image Segmentation Effect

As can be seen from figure 3-2, although the degree of adhesion of some of the verification code was very serious and even the existence of a common boundary conditions, but use of projection histogram can still be a good segmentation of the characters.

Since the neural network requires the input format to be consistent, it is necessary to make the size of segmented character image normalize the to $36 * 24$ pixels.

3.3 BP Neural Network

BP (Back Propagation) neural network is one of the most widely used neural networks which is called artificial neural network based on error back propagation algorithm. BP neural network can be adjusted according based on the propagation mode of error back propagation in order to get a multi-level feedback network, and have more less analysis in the input information and the predictability of output. BP neural network have a lot of advantages such as modeling method accurately, has fast transmission speed and has higher recognition efficiency, although the training of multi-layer network was difficult and the speed was slow, but the BP neural network could achieve good effect in character recognition. At the same time applied to character recognition, we need to extract the features of the characters to be identified and then use the obtained feature vector to train the neural network classifier. The extracted features directly affect the final classification performance which to a certain extent Limiting the final recognition rate.

3.3.1 Neural Network Topology

Neural network is a classifier with strong robustness which has the characteristics of fault tolerance, self-adaptability, parallel processing, its learning ability comes from its topological structure. Neural network topology has many different structures in which feed-forward network refers to the network in the input signal in one direction from one node to another node continuous transmission until reached the output layer. This experiment was used in single-layer feed-forward network topology. As shows Figure 4.

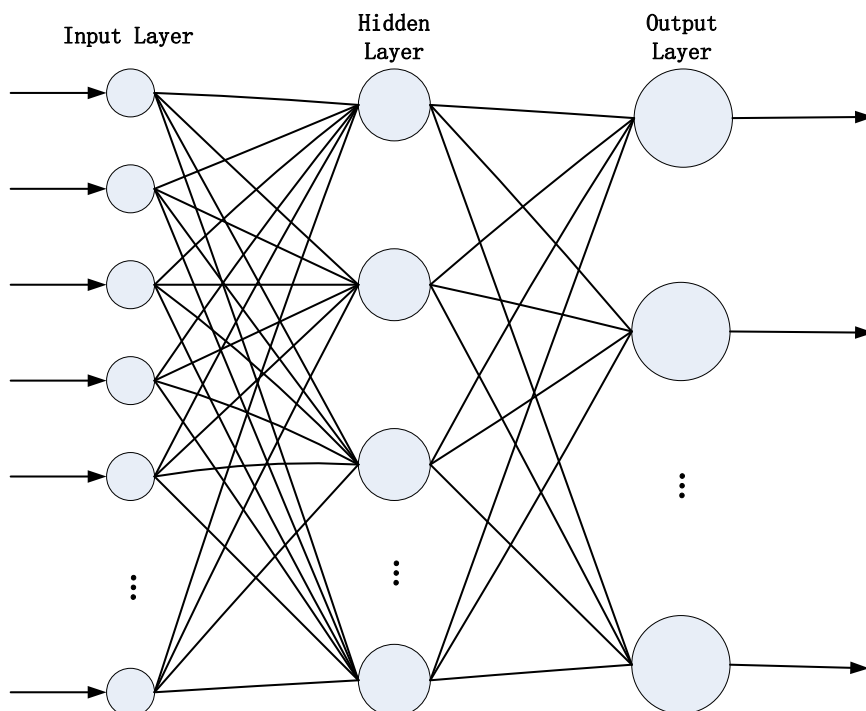


Figure 4. Single - Layer Feed-Forward Network Topology

Wherein the input layer inputs 864 pixel information about the single character image which is cut; the number of nodes in the hidden layer was 96 and many experiments show that the 96 hidden nodes can produce high correct rate, too many nodes would easily lead to over-fitting and slow down the training speed; the output layer consists of a linear neuron which has 37 target outputs included the 26 letter, 10

digits and an unknown (when the character can not recognition was marked unknown).

3.3.2 The Activation Function of Neural Network

Activation function refers to existence the nodes in the model contains many connectionstructure and each node represents a special output function. The activation function could transform the input signal into a single output signal that in order to facilitating the signal to propagate further over the network. It is the mechanism by which artificial neurons process information and transmit the information to the entire network. In this paper, the activation function used in the hidden layer was hyperbolic tangent function and in the output layer used the sigmoid function.

The hyperbolic tangent function (tanh) is the quotient of the hyperbolic sine function (sinh) and the hyperbolic cosine function (cosh). Formula (1) showsthe function formula:

$$\tanh(x) = \frac{1 - e^{-x}}{1 + e^{-x}} \quad (1)$$

The sigmoid function is very popular inANN which is a strictly increasing function andit exhibiting a good balance between linearity and nonlinearity. Formula (2) showsthe function formula:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (2)$$

4. Experimental Process and Result Analysis

4.1 Experimental Parameter Setting and Evaluation Index

The structure and activation function of the neural network were determined, then in the process of the experiment still need to configure various parameters. Table 1 shows parameters setting.

Table 1. Parameter Setting

Name	Initial Learning Rate	Momentum	Target Error	Maximum Number of Training
Value	0.0001	Learning Rate *0.9	0.002	256

The training type used in this experiment was random gradient method. After a number of experiments compared to the maximum number of training sets was set 256 had faster training speed and accuracy, but when the target error to the target error was the end of training, so the target error value should not be set too large.Training error trajectory and learning rate change shown in Figure 5.

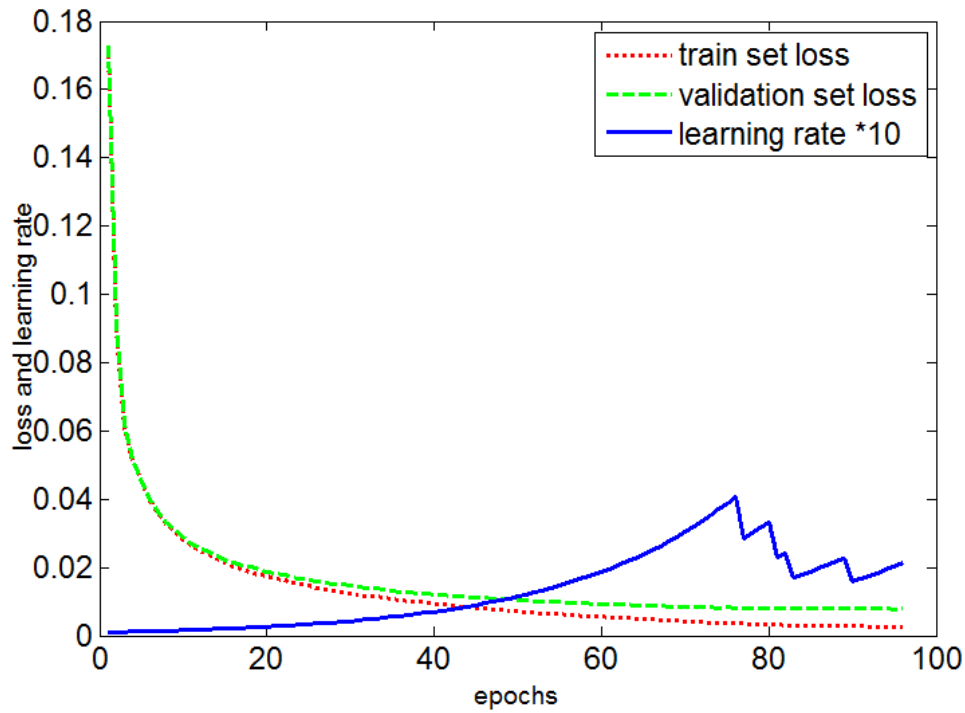


Figure 5. Training Error Trajectory

The evaluation index used in this paper was the correct rate of character recognition--Q which was the mean value of all the characters correct rate--P, P was show in Formula (3). TP was the number of characters which had correctly identified and the FN was the number of characters which wrong recognize to the other characters.

$$P = \frac{TP}{TP + FN} \quad (3)$$

4.2 Experimental Results and Analysis

There were 10,480 training samples in the experiment and the correct rate was increasing until the sample reaches a relatively stable state. The frequency of the characters contained in the sample was show in Table 2. We can see from the data in the table that apart 0, o, 1, I, had less frequency which characters were too hard to recognize for human and other characters had average proportion. Because the uneven sample will result in over-fitting training and make the fewer samples of characters has too low recognition accuracy rate to affect the overall accuracy.

Table 2. Character Frequency

Character	A	B	C	D	E	F	G
Frequency	2.86%	3.36%	3.36%	3.02%	2.96%	2.86%	3.03%
Character	H	I	J	K	L	M	N
Frequency	3.13%	0%	3.17%	3.28%	2.94%	2.43%	3.39%
Character	O	P	Q	R	S	T	U

Frequency	0.10%	3.18%	2.95%	3.29%	2.98%	2.99%	2.86%
Character	V	W	X	Y	Z	0	1
Frequency	2.9%	3.04%	3.32%	3.22%	3.29%	0.06%	0.15%
Character	2	3	4	5	6	7	8
Frequency	2.75%	2.94%	3.05%	3.27%	3.19%	2.97%	2.93%
Character	9						
Frequency	2.95%						

In order to obtain the accurate experimental results, we used 9-fold cross validation to divide all the samples randomly into 9 groups that mean one group was the test set in succession and the remaining 8 groups were the training set and the validation set. The final experimental results was used the average of 9 experiments. The accuracy of the nine experiments is shown in Table 3.

Table 3. Recognition Accuracy

Times	1	2	3	4	5	6	7	8	9
Q	0.931	0.921	0.916	0.917	0.915	0.921	0.924	0.918	0.926

From the above table calculated the correct rate of this experiment was 0.921. It was proved that this kind of character recognition algorithm based on BP neural network has a high recognition rate and low error rate which can meet the requirement of high precision, at the same time the method also has a certain degree of universality that for other types of CAPTCHA can achieve a good recognition effect and has certain practical value.

5. Conclusion

This paper introduces the CAPTCHA recognition algorithm of character adhesion based on BP neural network. The experimental results show that the proposed algorithm has high recognition rate through image preprocessing and character segmentation. The algorithm is simple to use, and it can remove the complex noise in the background well, and it can be better to segment the characters. It is also simple to train the neural network training classification, only a small amount of prior knowledge and good adaptability Advantages, for the protection of network security and verification code design and application of practical significance. But for the serious adhesion of the characters still exist in the case of poor division, recognition errors and so on. The speed of CAPTCHA innovation is accelerating, have more and more types, more and more complex content, not only have character adhesion and even has overlapping distorted characters. How to improve the accuracy of distorted sticky character recognition is the direction and focus for the next study step.

Acknowledgments

This research is supported by National Natural Science Foundation of China Under the Grant 61672210 and supported by the Henan Research Program of Foundation and Advanced Technology Under the Grant 162300410183.

References

- [1] S. Yongmei, H. Hua, "Research on Domain-independent Opinion Target Extraction", *International Journal of Hybrid Information Technology*, vol. 8, no. 1, (2015), pp. 237-248
- [2] L. von Ahn, M. Blum and J. Langford, "Telling Humans and Computer Apart Automatically", *Comm. Of the ACM*, vol. 46, (2003), pp. 57-60.
- [3] S. Murugavalli, SAK Jainulabudeen, Kumar GS, "Anuradha D. Enhancing security against hard AI problems in user authentication using CAPTCHA as graphical passwords", *International Journal of Advanced Research in Computer Science*, vol. 6, no. 24, (2016), pp.93-99.
- [4] B.B. Zhu, J. Yan, G. Bao, M. Yang, "Captcha as Graphical Passwords—A New Security Primitive Based on Hard AI Problems", *Information Forensics & Security IEEE Transactions*, (2014), vol. 9, no. 6, pp. 891-904.
- [5] A. Rusu, V. Govindaraju, "Handwritten CAPTCHA: using the difference in the abilities of humans and machines in reading handwritten words", *Computational Intelligence in Bioinformatics and Computational Biology*, 2004. CIBCB '04. Proceedings of the 2004 IEEE Symposium, (2004).
- [6] A.O. Thomas, V. Govindaraju, N.Y. Amherst, "Veneration and performance evaluation of synthetic handwritten CAPTCHAs", *Proc of the 1st Int Conf on Frontiers in Handwriting Recognition, ICFHR 2008*. Washington, DC: IEEE Computer Society, (2008).
- [7] R. Gossweiler, M. Kamvar, S. Baluja, "What's up CAPTCHA? A CAPTCHA based on image orientation", *Proc of WWW 2009*. New York: ACM, (2009).
- [8] H. Gao, H. Liu, D. Yao, X. Liu, U. Aickelin, "An Audio CAPTCHA to Distinguish Humans from Computers", *International Symposium on Electronic Commerce & Security*, (2010), pp. 265-269.
- [9] H. Gao, D. Yao, H. Liu, X. Liu, L. Wang, "A Novel Image Based CAPTCHA Using Jigsaw Puzzle", *15th International Conference on Computational Science and Engineering*, (2010).
- [10] M. Wang, T. Zhang, W. Jiang, H. Song, "The Recognition of CAPTCHA", *Journal of Computer & Communications*, (2014), vol. 2, no. 2, pp. 14-19.
- [11] S.Y. Huang, Y.K. Lee, G. Bell, Z. Ou, "A Projection-based Segmentation Algorithm for Breaking MSN and YAHOO CAPTCHAs", *Lecture Notes in Engineering and Computer Science*, (2008), vol. 2170, no. 1.
- [12] S.-Y. Huang, Y.-K. Lee, G. Bell, Z.H. Ou, "An efficient segmentation algorithm for CAPTCHAs with line cluttering and character warping", *Springer Science Business Media, LLC*, (2009), vol. 48, pp. 267-289.
- [13] J. Cui, L. Liu, G. Du, Y. Wang, "Skeletonization of Deformed CAPTCHAs Using Pixel Depth Approach", *Journal of Multimedia*, vol. 6, no. 6, (2012), pp. 526-533.
- [14] P. Guo, Y. W. Deng, H. Y. Zhang, "A captcha image recognition algorithm based on edit distance", *Key Engineering Materials*, (2011), 474-476: 2203-2207.
- [15] L. Zhang, S. G. Huang, Z. X. Shi, "CAPTCHA Recognition Method Based on RNN of LSTM", *Pattern Recognition and Artificial Intelligence*, vol. 24, no. 1, (2011), pp. 40-47.
- [16] L. Zhang, S. G. Huang, "A novel method to recognise closely connected captcha", *Advanced Materials Research*, (2012), 457-458: 620-627.
- [17] C. J. Chen, Y. Wei, W. P. Fang, "A study on captcha recognition", *Tenth International Conference on Intelligent Info*, (2014), pp. 395-398.
- [18] L. Yin, D. Yin, R. Zhang, "A Recognition Method for Distorted and Merged Text-Based CAPTCHA", *Pattern Recognition and Artificial Intelligence*, vol. 27, no. 3, (2014), pp. 235-241.
- [19] K. Chellapilla and P. Simard, "Using Machine Learning to Break Visual Human Interaction Proof", *Neural Information Processing Systems (NIPS)*, MIT Press, (2004).
- [20] G. Yin, L. Tao, "Verified Code Recognition Algorithm Based on SVM", *Computer Engineering and Applications*, (2011), vol. 47, no. 18, pp. 188-190.
- [21] L. Wang, R. Zahang, D. Yin, J. Zhan, C. Wu, "Breaking visual CAPTCHA of merged characters", *Computer Engineering and Applications*, vol. 47, no. 28, (2011), pp. 150-153.
- [22] H. DHuo, D. Li, "News Video Text Area Positioning Based on an Improved Trajkovic Corner Detector", *International Journal of Future Generation Communication and networking*, (2016).

Authors



Huo Hua, he is a professor with a Ph.D. in Information Engineering College and School of Software, Laboratory of Intelligent Computing & Application Technology for Big Data, Henan University of Science and Technology. His research interests in intelligent information processing and video semantic extraction.



Chang Guoqin, she is a postgraduate student of Laboratory of Intelligent Computing & Application Technology for Big Data, Henan University of Science and Technology. Her research interests in intelligent information processing and data mining.