

# Research on the Optimization of License Plate Character Segmentation

Yunwei Cai \*

School of Intelligent Manufacturing and Control Engineering Shanghai, Shanghai Second Polytechnic University, Shanghai, 200135, China

\* Corresponding Author Email: Yunweicai2022@126.com

**Abstract.** License plate character segmentation is an important step in license plate detection and recognition. With the development of machine learning technology, the segmentation algorithm of license plate characters based on clustering is also developed rapidly. However, the current clustering algorithm based on K-means does not consider the integrity of the characters and the horizontal difference between the characters on the license plate. This paper presents a weighted distance measurement method based on a connected domain. First, all pixels belonging to the same connected domain are naturally clustered into a class. Meanwhile, the horizontal and vertical distance measurements between the connected domains are scaled by the respective land weights, which makes the horizontal distance between the connected domains get more attention. This paper conducts extensive experiments on the collected data set of license plate characters, and the experimental results verify that compared with the K-means clustering algorithm, the weighted distance measurement method based on the connected domain has more accurate license plate character segmentation results.

**Keywords:** Vehicle License Plate Recognition; Character Segmentation; Weighted Distance Measurement.

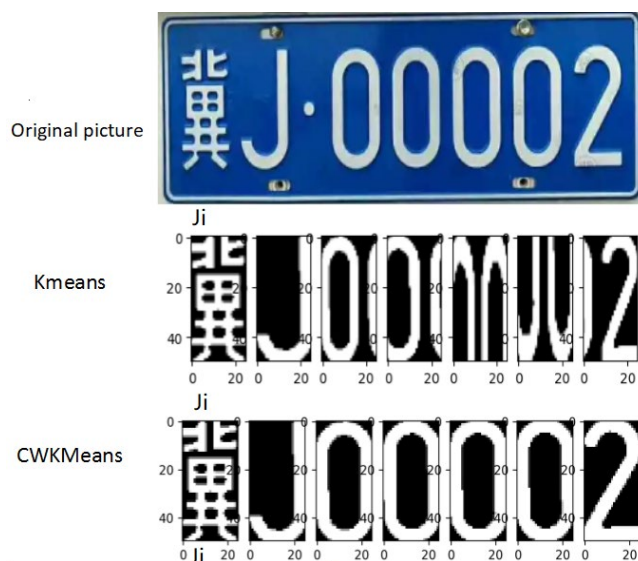
## 1. Introduction

Image processing technology plays an important role in the field of computer vision. As a traditional image processing task, license plate detection and recognition involve a variety of image processing technologies and machine learning algorithms. Nowadays, our car license plate generally consists of seven characters and a point, the height and width of the license plate characters are fixed, respectively 90mm and 45mm, the distance between seven characters is also fixed at 12mm, and the diameter of the point partition is 10mm [1]. In civil license plates, the arrangement of characters follows the following rules: the first character is usually the abbreviation of Chinese provinces and regions, expressed with Chinese characters; The second character is usually the code number of the issuing authority. The last five characters are a combination of letters and digits. The letters are a combination of twenty-four uppercase letters, and the digits are represented by digits ranging from 0 to 9. At present, Chinese license plates can be roughly divided into white characters on a blue background and black characters on yellow background. Special cars use black characters on a white background or white characters on black background, sometimes supplemented by red characters and so on. To simplify processing, only license plates with white characters on a blue background were considered in this study.

License plate character segmentation is an important part of license plate detection and recognition. In recent years, with the development of license plate detection and recognition, more novel license plate character segmentation algorithms appear. For example, Use sections, and then the horizontal segmentation method of the Hough transform fitting table is used to remove the influence of the top and bottom frame and rivet [2]. Based on the horizontal projection analysis of the vertical projection, the two steps of license plate character segmentation and recognition are taken as a whole statistical inference problem [3]. Through the R channel pretreatment method, other interference factors outside the license plate are eliminated, so that the license plate characters and background are completely separated [4]. Character height and interval characteristics obtained by preprocessing and connected domain are used to screen characters, so as to improve the anti-interference ability of the segmentation

algorithm to the border of license plates and fake license plates [5]. Based on CRF license plate character segmentation algorithm, license plate recognition is converted into license plate image column classification, and structured machine learning is carried out to solve the license plate character segmentation problem with low pixel quality [6]. The mean value outburst method is used to locate the license plate, and the K-means algorithm is used to cluster and segment the license plate characters, thus improving the accuracy of license plate character segmentation [7]. Through deep learning Faster R-CNN combined with K-means algorithm for license plate positioning, a CNN model for license plate character recognition is proposed [8]. But these methods do not consider character segmentation should pay more attention to the horizontal direction of the problem, also do not consider the overall character characteristics.

In morphology, letters and numbers in license plate characters form a connected domain, while Chinese characters generally contain connected domain parts of several radicals. In character segmentation, all pixels belonging to the same connected domain should be naturally divided into the same category. For example, multiply connected domains in Chinese characters should be clustered into one category. The traditional K-means clustering algorithm usually adopts the horizontal and vertical undifferentiated Euclidean distance (L2) or chessboard distance (L1), which leads to the horizontal distance and vertical distance between the clustered elements being treated equally, so that the upper and lower parts belonging to the same character may be clustered into different categories. However, the close horizontal distance causes parts belonging to different characters to be grouped together. As shown in Figure 1, the traditional K-means algorithm will not only divide the upper and lower parts of the same character (the “0” in the license plate), but also gather the two characters that are close to each other into a class (the horizontal distance between the “0” are close, which are aggregated into a class by K-means algorithm). This paper proposes a distance-weighted K-means clustering algorithm based on a connected domain (CWKMeans). First, all pixels belonging to the same connected domain are naturally clustered into a class. Meanwhile, the horizontal and vertical distance measurements between the connected domains are scaled respectively, which makes the horizontal distance between the connected domains get more attention. Specifically, the horizontal and vertical distance between two connected domains is respectively  $dist_x$ ,  $dist_y$ . Then the distance between the final connected domains is defined as  $dist = \alpha * dist_x + \beta * dist_y$ . Therefore, horizontal and vertical distances between the connected domains are scaled by different hyperparameters. In the experiment, we use the larger  $\alpha$  ( $\alpha=20$ ) and smaller  $\beta$  ( $\beta=1$ ) to make the clustering algorithm pay more attention to the horizontal distance of the connected domain.



**Figure 1.** Distance-weighted K-means clustering algorithm based on connected domain (CWKMeans) and traditional K-means algorithm for license plate character segmentation comparison of results

This paper collected the validity of the verification method of the license plate image data set. Through a large number of comparison experiments and ablation experiments, it is proved that the CWKMeans algorithm can cluster the same character into a class more accurately, and at the same time accurately segment adjacent characters.

## 2. The Method of Card Character Segmentation

### 2.1 License Plate Image Preprocessing

The preprocessing of license plate image includes gray level image, binarization and morphologic opening to remove part of the adhesive and independent noise points in the image. The pre-processing results are shown in Figure 2. OTSU algorithm is used for binarization and a 2X2 kernel is used for Open.

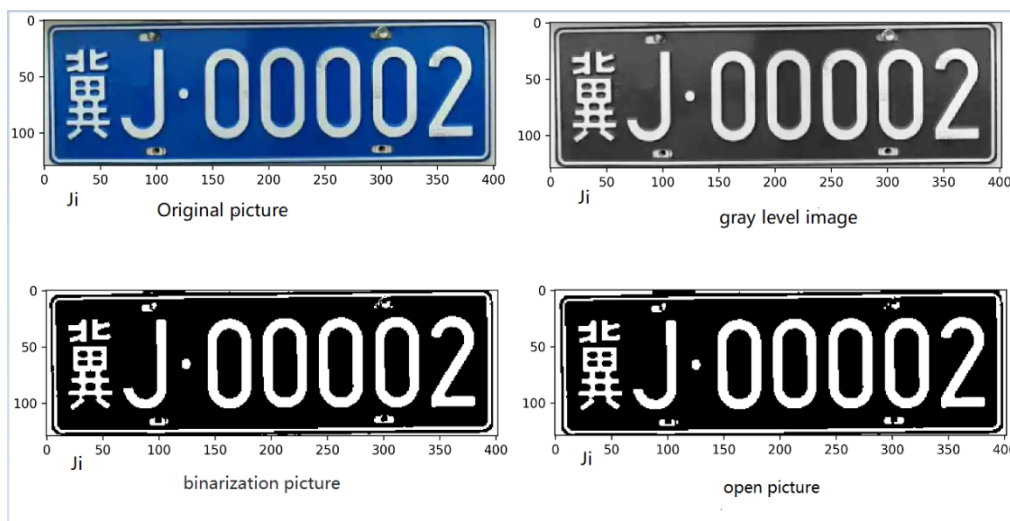


Figure 2. License plate image preprocessing

### 2.2 Remove the Plate Frame and Plate Alignment

Next, remove the border of the license plate picture (top, bottom, left, and right border) and the dot identifier in the license plate, and remove the top and bottom dot to align the license plate. First, remove the border of the license plate. Since the border of the license plate has connectivity, that is, the upper and lower border of the license plate and the left and right border are connected together, the connected domain whose width is greater than  $1/7$  of the width of the entire license plate is set as the border area in this paper. Second, remove the dots. Because there is a split dot between the second character and the third character in the license plate character, its width is small, and the width and height is the same. Let the coordinates of the center point of the connected domain be  $(x_{componet}, y_{componet})$ . Where,  $x_{componet}$  is the horizontal central coordinate point of the connected domain, and  $y_{componet}$  is the vertical central coordinate point of the connected domain. The width and height of the license plate are  $I_{width}$  and  $I_{height}$  respectively. In this paper, the horizontal and vertical coordinate error of the center point of the connected domain is set at 2 pixels, namely  $|x_{componet} - y_{componet}| < 2$ , the horizontal center coordinate point is in the range of  $2/7$ - $4/7$  horizontal of the license plate, namely  $\frac{2}{7} < \frac{x_{componet}}{I_{width}} < \frac{4}{7}$  and the vertical center point coordinates are in the range of  $2/5$ - $3/5$  vertical of the license plate, namely  $\frac{2}{5} < \frac{y_{componet}}{I_{height}} < \frac{3}{5}$ . Then, license plate alignment is performed. In this paper, the line projection method is adopted for license plate alignment [9]. The sum of row pixel values of the license plate was counted and denoted as the vector  $v_{row}$ . The minimum value point of  $v_{row}$  and the mean value of the mean point of  $v_{row}$  were set as the threshold. The  $v_{row}$  was scanned, and the threshold value was used as the boundary to divide the

cut segment, and the largest cut segment was used as the license plate alignment result. The result of preprocessing the license plate is shown in Figure 3.

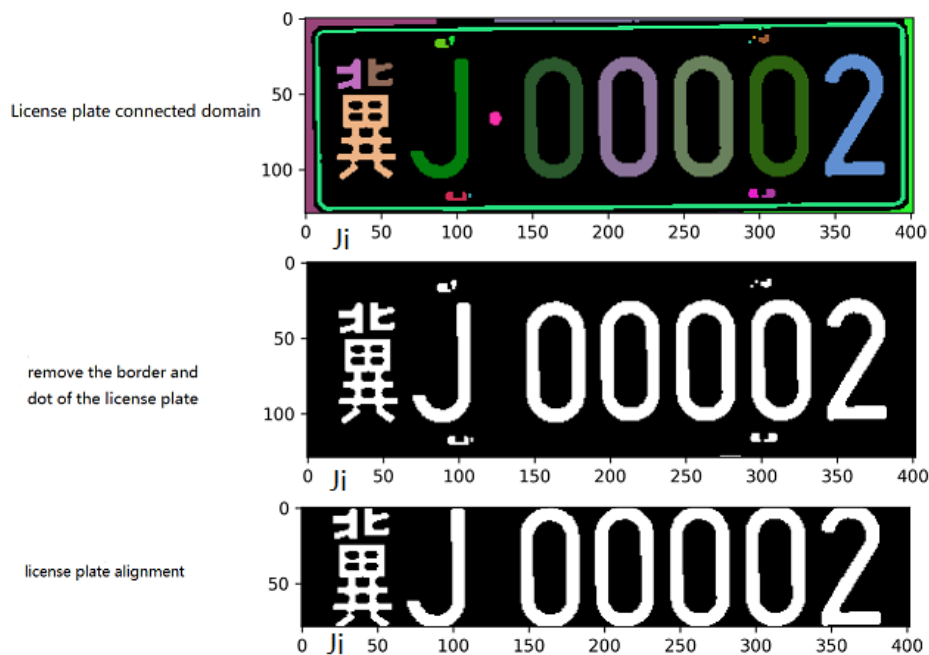


Figure 3. Remove plate border and plate alignment

### 2.3 Character Clustering based on K-means

Based on Section 2 above, K-means character clustering is carried out for the aligned license plate binary image. The specific clustering process is as follows:

- The alignment of the pixels in binary image  $I_{align}$  represented as pixels to set the DataSet =  $\{(x_i, y_i) | I_{align}(y_i, x_i) = 255\}, i = 1, 2, \dots, N$ , where  $N$  represents the number of non-zero pixel points;
- The clustering centers were initialized and the number of clustering centers was  $k$ . The locations of  $k$  clustering centers were in the longitudinal center of the license plate and evenly distributed in the horizontal direction. The  $k$  clustering centers were represented as  $c = \{c_1, c_2, \dots, c_k\}$ ;
- For each pixel point  $p_i = (x_i, y_i), i = 1, 2, \dots, N$  in DataSet dataset, the distance from it to  $k$  clustering centers was calculated and assigned to the corresponding class of the cluster center with the smallest distance;
- For each category  $c_i (i = 1, 2, \dots, k)$ , to calculate its clustering center  $c_i = \frac{1}{|c_i|} \sum_{p_i \in c_i} p_i$  (namely belongs to the class of all samples of the center of mass);
- Repeat the above steps c) and d) until the clustering center does not change.

### 2.4 Segmentation of License Plate Characters based on CWKMeans Algorithm

#### 2.4.1 Connected Domain Statistics

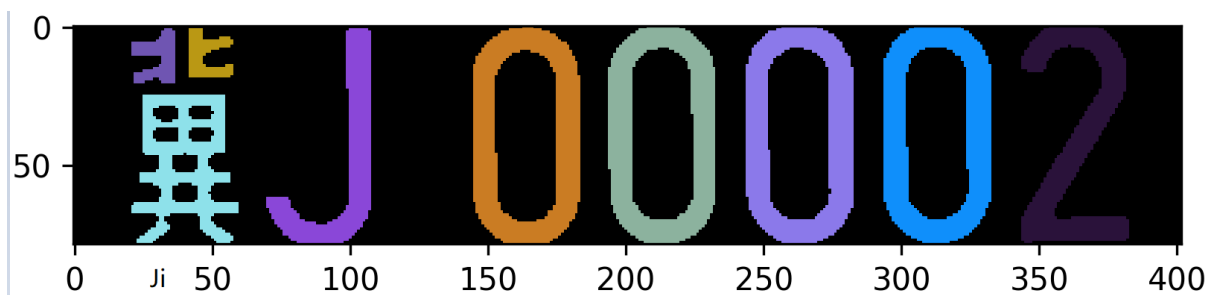


Figure 4. Connected domain statistics

All connected domains are obtained through the 8-neighborhood domain, denoted as  $T = \{T_1, T_2, \dots, T_M\}$ , where  $M$  represents the number of connected domains, and any two pixels in each connected domain  $T_j, j = 1, 2, \dots, M$  can finally find each other through the pixels in the 8 domains. In this paper, we took the set of central points of each connected domain as the cluster point set, that is,  $DataSet = \{center(T_1), center(T_2), \dots, center(T_M)\}$ . As shown in Figure 4.

### 2.4.2 Initialize the Cluster Center

Seven clustering centers were horizontally and vertically, uniformly set  $c = \{c_1, c_2, \dots, c_k\}, k=7$ , same as flow b) in section 3. As shown in Figure 5.

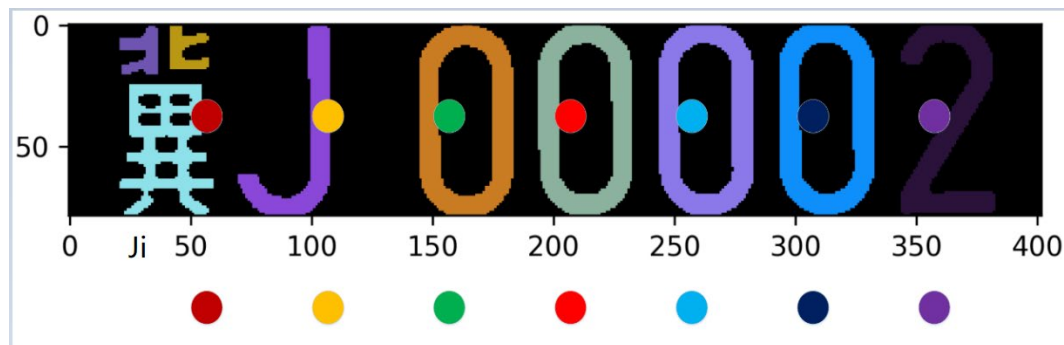


Figure 5. Initialize the cluster center point

### 2.4.3 Define the Weighted Distance Measure

The distance between any two two-dimensional coordinate points  $p_1 = (x_1, y_1), p_2 = (x_2, y_2)$  is expressed as  $d(p_1, p_2) = \alpha|x_1 - x_2| + \beta|y_1 - y_2|$ . Where  $x$  is the horizontal coordinate in the image,  $y$  is the vertical coordinate in the image, and  $\alpha$  and  $\beta$  are the hyperparameters, in this paper, we set  $\alpha$  as 20,  $\beta$  as 1. Different proportion parameters will make the clustering between connected domains pay more attention to the horizontal distance so that different characters are separated and the connected domains of the same character are more aggregated (different connected domains of the same character show small horizontal distance and large vertical distance).

### 2.4.4 CWKMeans Clustering Iterative Process

Based on Section 2 above, K-means character clustering is carried out for the aligned license plate binary image. The specific clustering process is as follows:

- The alignment of the pixels in binary image  $I_{align}$  represented as pixels to set the  $DataSet = \{(x_i, y_i) | I_{align}(y_i, x_i) = 255\}, i = 1, 2, \dots, N$ , where  $N$  represents the number of non-zero pixel points;
- The clustering centers were initialized and the number of clustering centers was set as  $k$ . The locations of  $k$  clustering centers were in the longitudinal center of the license plate and evenly distributed in the horizontal direction. The  $k$  clustering centers were represented as  $c = \{c_1, c_2, \dots, c_k\}$ ;
- For each pixel point  $p_i = (x_i, y_i), i = 1, 2, \dots, N$  in DataSet dataset, the distance from it to  $k$  clustering centers was calculated and assigned to the corresponding class of the cluster center with the smallest distance;
- For each category  $c_i (i = 1, 2, \dots, k)$ , to calculate its clustering center  $c_i = \frac{1}{|c_i|} \sum_{p_i \in c_i} p_i$  (namely belongs to the class of all samples of the center of mass);
- Repeat the above steps c) and d) until the clustering center does not change.

## 3. The Experiment

### 3.1 Experimental Setup

Based on Python language environment, using Numpy, Opencv, matplotlib and other libraries to achieve image operation and visualization; The Python IDE uses Pycharm software and runs in the

python environment created by Anaconda, combining K-means clustering algorithm and CWKMeans clustering algorithm to achieve license plate character segmentation. In this paper, 100 license plate images with white letters on blue background were collected and saved in img format as data sets, including numbers and letters of multiple provinces with high coverage, which enhanced the universality and reliability of the experimental results.

### 3.2 Analysis of Visual Experimental Results



Figure 6. Comparison of operation results of two algorithms

As shown in Figure 6, compared with the traditional K-means-based clustering algorithm, the distance-weighted K-means clustering algorithm proposed in this paper can scale the horizontal and vertical distance measurements between the connected domains respectively, so as to avoid the two characters with close horizontal distance being clustered into one category, and prevent the single character with wide distance from being divided into two categories. As shown in Figure 7, the length and width of the characters processed by the distance-weighted K-means clustering algorithm based on the connected domain are more regular, and there will not be a large area gap between different clustering modules so that the clustering is more beautiful, the character thickness is more uniform, and the original appearance of the license plate characters can be more accurately restored.

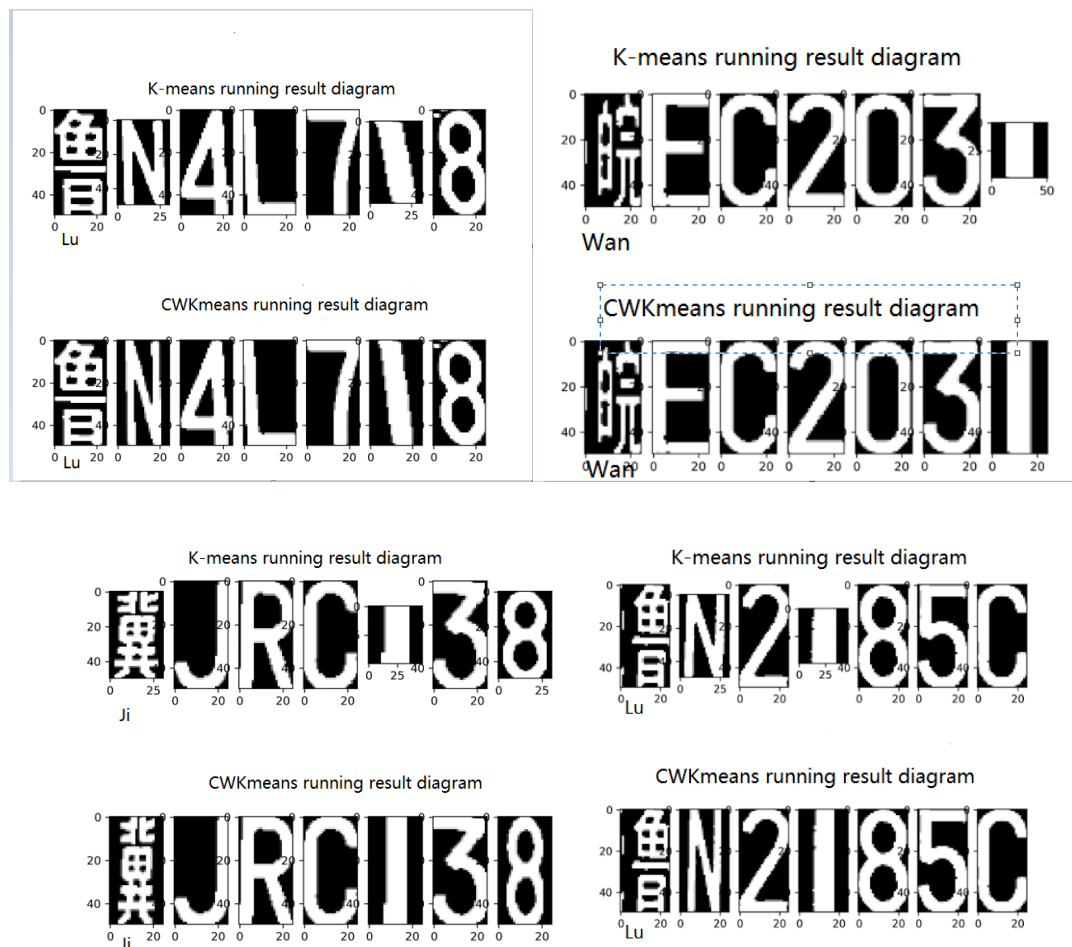


Figure 7. Comparison of character segmentation effect

#### 4. Summary and Outlook

In this paper, the unweighted K-means method is compared with the randomly initialized clustering center. After the weighted K-means method and self-selection of more appropriate initial clustering, the number of iterations is reduced and the efficiency of the final clustering results is improved. Although compared with the traditional K-means clustering algorithm, this experiment is more accurate in segmentation and more regular in word size, perfect correction or recovery cannot be achieved when characters are tilted or damaged and folded. In the following work, the author will try to introduce and write library functions and codes that can adjust character angles and recover missing parts.

#### References

- [1] Gao Yong, Character Segmentation and Recognition in License Plate Recognition System, Anhui University, 2007.
- [2] ZHANG Yungang, License Plate Character Segmentation Algorithm based on Hough transform and Prior Knowledge. Beijing: Institute of Information Processing, Department of Automation, Tsinghua University, 2004.
- [3] JIAO Hui-hua. Research on Character Segmentation of License plate image Based on Vertical Projection Segmentation Method. Haikou: Qiongtai Normal University, 2022.
- [4] Lai Daoliang. Research on Image Preprocessing method for License Plate Character Segmentation. Chengdu: School of Economic Information Engineering, Southwestern University of Finance and Economics, 2018.

- [5] Shi Longzhao. Character Segmentation Algorithm of Complex License Plate Based on Connected domain. Fuzhou: College of Physics and Information Engineering, Fuzhou University,2016.
- [6] Fu J Q. Low image quality License plate character segmentation based on Conditional Random field. Shanghai: Institute of Media Computing, School of Computer Science and Technology, Fudan University, 2014.
- [7] Li Jin. Improved K-means Segmentation of License Plate Characters. Nanjing. School of Electronic and Information Engineering, Nanjing University of Technology,2015.
- [8] LI Xiangpeng. License Plate Location and Recognition Method Based on Deep Learning [J]. Nanchang. School of Information Engineering, Nanchang University, 2019.
- [9] WU Genxing. Using Projection Method to Cut two-dimensional Line Segment Rectangular Window. Hangzhou. College of Mechanical and Electrical Engineering, China Jiliang University,2007.