

Traffic Sign Recognition Based on CNN and Twin Support Vector Machine Hybrid Model

Yang Sun, Longwei Chen

College of Statistics and Mathematics, Yunnan University of Finance and Economics, Kunming, China

Email: sunyangxyr@163.com, ZZ1237@ynufe.edu.cn

How to cite this paper: Sun, Y. and Chen, L.W. (2021) Traffic Sign Recognition Based on CNN and Twin Support Vector Machine Hybrid Model. *Journal of Applied Mathematics and Physics*, 9, 3122-3142. <https://doi.org/10.4236/jamp.2021.912204>

Received: November 25, 2021

Accepted: December 21, 2021

Published: December 24, 2021

Copyright © 2021 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

With the progress of deep learning research, convolutional neural networks have become the most important method in feature extraction. How to effectively classify and recognize the extracted features will directly affect the performance of the entire network. Traditional processing methods include classification models such as fully connected network models and support vector machines. In order to solve the problem that the traditional convolutional neural network is prone to over-fitting for the classification of small samples, a CNN-TWSVM hybrid model was proposed by fusing the twin support vector machine (TWSVM) with higher computational efficiency as the CNN classifier, and it was applied to the traffic sign recognition task. In order to improve the generalization ability of the model, the wavelet kernel function is introduced to deal with the nonlinear classification task. The method uses the network initialized from the ImageNet dataset to fine-tune the specific domain and intercept the inner layer of the network to extract the high abstract features of the traffic sign image. Finally, the TWSVM based on wavelet kernel function is used to identify the traffic signs, so as to effectively solve the over-fitting problem of traffic signs classification. On GTSRB and BELGIUMTS datasets, the validity and generalization ability of the improved model is verified by comparing with different kernel functions and different SVM classifiers.

Keywords

CNN, Twin Support Vector Machine, Wavelet Kernel Function, Traffic Sign Recognition, Transfer Learning

1. Introduction

The rapid and accurate identification of traffic signs plays a vital role in auto-

matic-assisted driving, and also plays an important role in improving safe driving. In the actual outdoor automatic driving, there are many unstable factors that will increase the difficulty of the system to identify the traffic scene. Due to the influence of the rapid movement of vehicles, such as image blurring, image fading, weather and light environment, the traffic sign recognition becomes more complex compared with the still life recognition, and at the same time, it puts forward higher requirements for the recognition model.

To solve the problem of traffic sign recognition, researchers around the world have proposed a large number of algorithms, which are mainly divided into two categories: traditional algorithms and deep learning algorithms. Traditional algorithms mainly rely on manual extraction of image features according to image morphology, and combine with machine learning algorithm for recognition. Image features mainly include HSV color space, directional gradient histogram, local binary mode and Gabor feature, etc. In feature extraction methods, there are mainly SIFT transform, Hough transform and Wavelet transform. For the extracted features, SVM, Adaboost and other classifiers are used to complete the identification of traffic signs. Because the acquisition of traffic sign images is accompanied by a complex actual environment, the design of feature map and feature vector becomes complicated and does not have good versatility, which means that the traditional manual extraction of feature map has been unable to meet the growing requirements. In comparison, the recognition system based on deep learning is a method with better effect and more applications. Among them, the convolutional neural network and its improved network have broken many records in many image classification competitions. It avoids the difficulty of manual design of feature extractor and has excellent performance in the feature extraction of traffic signs.

The powerful classification ability of CNN comes from its ability to learn the features of samples from a large number of samples. The whole network expresses the mapping relationship between original images and features. This also shows that CNN is a powerful feature extractor. According to different problems, different learning is carried out to get different network parameters. Numerous studies and experiments have shown that the network performance can be significantly improved by deepening the number of layers and using smaller convolution cores. VGGNet-16, VGGNet-19 and GoogLeNet are typical examples. But this is also accompanied by longer training times, which greatly limits the use of the network. In order to overcome the problems of insufficient training samples and too long training time, many pieces of literature migrate classifiers trained by large similar data sets to target problems and achieve good results. Transfer learning is also widely used in image recognition and natural language processing. Many studies have shown that using support vector machine to replace the softmax classifier at the last layer of the network can produce a better classification effect, and this method has been widely used in many fields.

Inspired by this, this paper uses the transfer learning method to fine-tune the structure of VGGNET-16 network, and designs a convolutional neural network model to extract traffic sign features. Aiming at the traditional CNN-SVM model, a new CNN-WTWSVM traffic sign recognition model was constructed by introducing the twin support vector machine (SVM) based on wavelet kernel function as a classifier for the first time in order to accelerate the calculation speed and improve the calculation accuracy. At the same time, the characteristics of network output are normalized to reduce over-fitting and achieve higher precision recognition.

2. Related Works

The traditional traffic sign recognition method mainly extracts the features according to the image and then analyzes the extracted features. This method is applicable to a small range and the operation is complicated. The research on this method has achieved fruitful results. In Literature [1], a multi-layer perceptron based on color threshold detection and neural network recognition is proposed to realize the recognition of traffic signs. Zhu Shuangdong *et al.* [2], put forward the concept of color shape pairs and built the color-geometric model of traffic signs to realize the detection and classification of traffic signs. In Literature [3], traffic signal detection, tracking and recognition system based on Adaboost training and Bayesian classification is proposed. The framework is based on the training of Adaboost to obtain a set of Haar wavelets for signal detection, and the Gaussian probability model is used to identify signals. Miao Xiaodong [4] *et al.* improved the multi-scale Log-Gabor wavelet for feature extraction, and finally realized classification by using the optimized SVM classifier. The experiment showed that the effect was good. Xu Shaoqiu *et al.* [5] simplified and decomposed the outer edge of the color segmentation output through discrete curve evolution, and then compared the smallest geometric difference with the template to identify the shape of the sign. This algorithm can realize reliable shape recognition in complex traffic scenes. Literature [6] proposes a shape-oriented algorithm, which shows good robustness for rotation and deformation in the process of image acquisition. In Literature [7], CIECAM97 model was used to segment and classify color regions, and FOSTS model was used to extract shape features, thus constructing a recognition framework with good performance in still images. Literature [8] proposed a three-step road sign recognition system based on color segmentation, shape recognition and neural network classification. Sun Guangmin *et al.* [9] built a traffic sign recognition system under natural conditions based on the color and geometric attributes of traffic signs themselves. Literature [10] used neural network to identify traffic sign data sets in the United States and Europe only based on shape detection, and achieved 90% accuracy on both data sets. Shams *et al.* [11] proposed a multi-class traffic sign recognition system based on BoW (Bag of Word) model. However, neglecting the spatial information is the weakness of the model, and then adding the spatial

histogram to retain the required spatial information. Finally, the test effect is generally good. At the same time, HOG (Histogram of Oriented Gradients) [12] and SIFT (Scale Invariant Feature Transform) [13] are also widely used in the identification task of traffic signs. In Literature [14], SIFI features are used to detect and describe the key points of the invariance of image scale and rotation, and then support vector machine is used to complete classification. A traffic sign recognition system based on SIFI and SURF [15] was established by the method in literature [16]. Multi-layer perceptrons were used for classification, and excellent classification effect was achieved. There are many other methods of traffic sign recognition based on HOG, such as literature [17] [18] [19] [20]. Abedin *et al.* [21] constructed an artificial neural network classifier based on HOG feature and accelerated robustness feature.

Compared with the traditional manual feature extraction classification method, deep learning is widely used in traffic sign recognition and has a better effect. Among them, convolutional neural network has a good performance in feature extraction. This method was first won many times in various computer vision challenges [22] [23], and later was applied in various fields of image processing. Literature [24] proposed a network framework with crossing layers and supplemented the framework with grayscale images. Finally, the accuracy was improved to 99.17%. In the work in Reference [25], the author proposed a method using local response normalization and HLSGD to train the Convolutional Neural Network, and the accuracy was 99.65%. Yin Shihao *et al.* [26] built a new convolutional neural network for traffic sign recognition by adding residual network, which improved the accuracy to 99.67%. Hu Wenzheng *et al.* [27] introduced branch convolutional neural network into deep convolutional neural network, and their experiments showed that traffic sign recognition could be carried out in a relatively shallow neural network. Zeng *et al.* [28] used the full connection layer in the convolutional neural network as a classifier to replace the extreme learning machine [29], and the accuracy reached 99.40%. In the early stage, Mattias *et al.* [30] developed a feature dimensionality reduction iterative nearest neighbor classifier based on iterative nearest neighbor linear prediction, with an accuracy of 98.53% on GTSRB and 98.32% on BTSC. Literature [31] proposed a two-module framework in which the maximum color probability extremum region was extracted as the recognition mark, and then the SVM classifier was trained. In the classification module, CNN detection symbols were used to classify the subclasses of each super class. This method greatly improved the computational efficiency.

Lian *et al.*'s work is also worth mentioning in traffic sign recognition. Inspired by speech recognition, Literature [32] proposed a frequency-selective assisted CNN model based on frequency domain. Jang *et al.* [33] used a spatial varying-voltage network to improve the geometric invariance of traffic sign symbols. Deng Zhidong *et al.* [34] proposed a relatively comprehensive method for detection and recognition of plane objects, which could achieve an average run-

ning time of 33.25 ms per frame for traffic signs within 100 meters. In Literature [35], the author proposed an overall framework based on SVM, which was used for the classification after convolutional neural network. Literature [36] proposed a weakly supervised measurement learning (WSMLR) method based on latent structural support vector machine (SVM). Chen *et al.* [37] proposed a road sign feature extraction method based on Gaussian-Hermite invariant moment (GHIMS). Then, the GHIMS features of different orders are transferred to BP neural network as vectors to realize the recognition of traffic signs.

Recent research on traffic sign recognition has achieved higher computational efficiency and recognition accuracy. Lu, Wang *et al.* [38] proposed an embedded multi-tasking learning model with multi-pattern tree structure called M2-TMTL to select visual features between and within patterns. In this method, two structured sparse-induced specifications are introduced into the least square regression. One of the specifications can be used not only to select the mode of a feature, but also to select features within the mode. The model is solved effectively by using the Multiplier Alternating Direction Method (ADMM). A large number of experiments on common benchmark datasets show that the algorithm has the same performance as several state-of-the-art methods, but with less computational and memory costs. Arcos-Garcia *et al.* [39] proposed a two-stage deep neural network including convolution and spatial converter network layer for image classification, with an accuracy of 99.71% on the GTSRB dataset. Zeng, Xu *et al.* [28] applied the extreme learning machine of depth perception features to deep learning networks and proposed a new traffic symbol recognition method called DP-KELM, which achieved 99.54% accuracy on GTSRB data sets. In the work of Zhang *et al.* [40], two new lightweight networks are proposed, which can achieve high identification accuracy while retaining fewer trainable parameters. Knowledge distillation transfers knowledge from a training model called a teacher network to a smaller model called a student network. On the GTSRB and BTSC traffic sign datasets, the identification accuracy of 99.61% and 99.13% can still be achieved by using smaller parameter sizes.

3. Proposed Methodology

3.1. Convolutional Neural Networks

Convolutional neural network is widely used in image processing and belongs to deep feed-forward learning artificial neural network. Convolutional neural network carries out a series of convolution and compression operations on the input image to obtain the feature map of higher abstract level, which can better extract the features of the image. Its powerful feature represents the learning ability, and it is invariant to the scale scaling, translation, rotation and other forms of the image, which makes the network, have a good learning ability. Another advantage of convolutional neural network is that it can effectively reduce the individual differences and obvious data noise of the data. It does not need to preprocess the image too much, and the data can be directly input into

the network for feature extraction.

Convolutional neural network is generally composed of convolutional layer, pooling layer, local response normalization layer, full connection layer and Soft-max layer. The convolutional layer carries out convolutional filtering operation on the input image through the coil nodule to extract its image features, and the result forms the feature map of this layer through activation function. Feature map not only reflects the size of the feature value, but also reflects its relative position relationship. For image feature extraction, it can retain more useful feature information. The convolution operation is usually accompanied by the operation of the convolution surface. The calculation formula of the convolution surface is as follows:

$$Y_i^l = f\left(\sum_{i=1}^D x_i^l * w_i^l + b_i^l\right) \tag{1}$$

where, l represents the number of layers, i represents the number of output convolution layers in this layer, and $f(\cdot)$ represents the activation function of this layer. $x_i^l * w_i^l$ represents the calculation of the i -th convolution surface, and the single convolution calculation formula is:

$$x * w = \sum_{s=1}^m \sum_{t=1}^n x_{i+m-s, j+n-t} \cdot w_{st}, \quad 1 \leq i \leq M - m + 1, 1 \leq j \leq N - n + 1 \tag{2}$$

where, $x = M \times N$, $w = m \times n$, $M \geq m$, $N \geq n$.

There are many kinds of activation functions in convolutional neural network. In image processing, the correction linear unit ReLU is generally selected, and its formula is as follows:

$$f = \max(0, x) \tag{3}$$

The pooling layer, namely the lower sampling layer, compacts the data obtained in the convolutional layer to simplify the computational complexity of the network. At the same time, the compression of feature images will lose part of the data and affect the accuracy of the network, so it is necessary to make up for this by deepening the network. Generally there are two types of average subsampling layer and maximum subsampling layer. The maximum pooling operation is to take the maximum value of the feature value in the neighborhood, which can retain more image texture information. In this paper, the method of maximum pooling is adopted. Its formula is as follows:

$$\text{maxdown}(G_{\lambda, \tau}^x(i, j)) = \max\{x_{st}, (i-1) \cdot \lambda + 1 \leq s \leq i \cdot \lambda, (j-1) \cdot \tau + 1 \leq t \leq j \cdot \tau\}$$

where, $G_{\lambda, \tau}^x(i, j) = (x_{st})_{\lambda \times \tau}$, $\lambda \times \tau$ is the size of the chunk.

Local response normalization is a method introduced to improve network performance and reduce network over-fitting. The local response normalization is calculated by the values of several adjacent convolution surfaces at position (x, y) , and the formula is:

$$y_{m,n}^l = x_{m,n}^l / \left(k + \alpha \sum_{j=\max(0, l-n/2)}^{\min(N-1, l+n/2)} (x_{m,n}^j)^2 \right)^\beta \tag{4}$$

where, N is the total number of convolutional surfaces, n is the number of adjacent surfaces, k, α, β is the tunable parameter, and $x_{m,n}^l$ represents the value at the position of (m, n) on the l convolutional surface.

The full-connection layer is a network layer that maps the distributed features learned from the convolution layer to the sample label space. Its essence is to transform from one feature space to another feature space. In the convolutional neural network, the full connection layer mostly appears in the last several layers to integrate the local information with category distinction between the convolutional layer and the pooling layer, and complete the feature weighting. Meanwhile, in the transfer learning task, the full connection layer can guarantee the transfer of the feature representation ability of the convolutional layer network. The Soft-max layer is a special kind of full connection layer that maps the output of multiple neurons between (0, 1), like the probability of the output of a single neuron, for the final implementation of classification.

With the continuous advancement of deep learning research, new models such as VGGNet and GoogLeNet appear successively, which verify that deep network can indeed improve the performance of network. The CNN network used in this paper takes VGG16 network model as its basic structure. The main idea of the network structure is to increase the depth of the network and reduce the size of the convolution kernel, so as to extract the image features more carefully. The VGG16 network in this paper has been trained on the ImageNet dataset in the source domain, and the network super parameters are obtained through transfer learning. Then the output layer structure of the network is fine-tuned to make the output image features migrate to the target domain. The network diagram is shown in **Figure 1**.

3.2. The Nonlinear Twin Support Vector Machine

Proposed by Vapnik in 1995, support vector machine (SVM) is a machine learning method based on statistical learning theory. Its advantages are shown in solving problems of small samples, nonlinear and high-dimensional pattern recognition. Based on the principle of structural risk minimization, this method has better generalization ability. The traditional support vector machine uses a hyper-plane

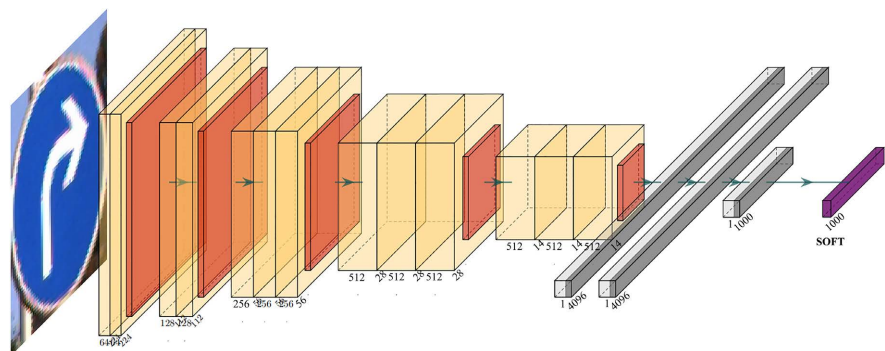


Figure 1. Vgg16 network structure diagram.

to divide two kinds of samples, and requires that the distance between the two kinds of samples formed by the sample points closest to the hyperplane reach the maximum. In order to improve the training speed of traditional support vector machines, Jayadeva *et al.* proposed a new machine learning algorithm called twin support vector machines (SVM) in 2007. The difference is that twin support vector machines try to find two non-parallel hyperplanes to divide the data set, and require one hyperplane to be as far away from one kind of sample as possible, and at the same time to be as close to the other kind of sample points as possible. It can be said that twin support vector machines are two quadratic programming problems similar to support vector machines, and theoretically, the training speed of twin support vector machines should be 4 times that of the original method.

For the training set $M \times N$, it contains two types of training samples, namely m_1 positive training samples in N -dimension, expressed as $m_1 \times N$, and m_2 negative training samples in N -dimension, expressed as $m_2 \times N$. The purpose of TWSVM is to find two non-parallel hyperplanes in the space of $M N$ -dimensional samples, which can not only correctly classify two different kinds of training samples, but also make each hyperplane as close to one kind of sample points as possible and as far away from the other kind of sample points as possible. The two hyperplanes are expressed as:

$$x^T w_1 + b_1 = 0, \quad x^T w_2 + b_2 = 0 \tag{5}$$

Since the training sample space studied in this paper is nonlinear, for the non-linear case, the kernel function is introduced in the same way as the traditional SVM, and the hyperplane based on the kernel space can be expressed as:

$$K(x^T, C^T)u_1 + b_1 = 0, \quad K(x^T, C^T)u_2 + b_2 = 0 \tag{6}$$

Then, the twin SVM in the partial linear case can be reduced to the following two quadratic programming problems:

$$\begin{aligned} \min_{w^{(1)}, b^{(1)}, \xi^{(2)}} & \frac{1}{2} \|K(A, C^T)w^{(1)} + e_1 b^{(1)}\|^2 + c_1 e_2^T \xi^{(2)} \\ \text{s.t.} & -(K(B, C^T)w^{(1)} + e_2 b^{(1)}) \geq e_2 - \xi^{(2)}, \xi^{(2)} \geq 0 \end{aligned} \tag{7}$$

$$\begin{aligned} \min_{w^{(2)}, b^{(2)}, \xi^{(1)}} & \frac{1}{2} \|K(B, C^T)w^{(2)} + e_2 b^{(2)}\|^2 + c_2 e_1^T \xi^{(1)} \\ \text{s.t.} & (K(A, C^T)w^{(2)} + e_1 b^{(2)}) \geq e_1 - \xi^{(1)}, \xi^{(1)} \geq 0 \end{aligned} \tag{8}$$

where c_1 and c_2 are penalty coefficients; e_1, e_2 represent two column vectors that are all 1's; $A = [x_1^{(1)}, x_2^{(1)}, \dots, x_{m_1}^{(1)}]^T$, $B = [x_1^{(2)}, x_2^{(2)}, \dots, x_{m_2}^{(2)}]^T$, $x_j^{(i)}$ represents the j -th sample of class i . $C = [A^T, B^T]^T$, A and B represent two types of training samples respectively.

Like the traditional SVM, the original problem is transformed into a dual problem for solving. First, introduce the Lagrange multiplier α and β , then

$$\begin{aligned}
 &L(w^{(1)}, b^{(1)2}, \xi^{(2)}, \alpha, \beta) \\
 &= \frac{1}{2} \left(K(A, C^T) w^{(1)} + e_1 b^{(1)} \right)^T \left(K(A, C^T) w^{(1)} + e_1 b^{(1)} \right) \\
 &\quad + c_1 e_2^T \xi^{(2)} - \alpha^T \left(- \left(K(B, C^T) w^{(1)} + e_2 b^{(1)} \right) + \xi^{(2)} - e_2 \right) - \beta^T \xi^{(2)}
 \end{aligned} \tag{9}$$

$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_{m_2})^T$, $\beta = (\beta_1, \beta_2, \dots, \beta_{m_2})^T$, with the KKT condition, the original problem can be transformed into a dual problem:

$$\begin{aligned}
 \max_{\alpha} \quad & e_2^T \alpha - \frac{1}{2} \alpha^T R (S^T S)^{-1} R^T \alpha \\
 \text{s.t.} \quad & 0 \leq \alpha \leq c_1
 \end{aligned} \tag{10}$$

where $S = \begin{bmatrix} K(A, C^T) & e_1 \end{bmatrix}$, $R = \begin{bmatrix} K(B, C^T) & e_2 \end{bmatrix}$, let $z_1 = \begin{bmatrix} w^{(1)} & b^{(1)} \end{bmatrix}^T$, then $z_1 = -(S^T S)^{-1} R^T \alpha$.

Similarly, after the same change for the second problem, its dual problem can be obtained as follows:

$$\begin{aligned}
 \max_{\gamma} \quad & e_1^T \beta - \frac{1}{2} \beta^T L (N^T N)^{-1} L^T \beta \\
 \text{s.t.} \quad & 0 \leq \beta \leq c_2
 \end{aligned} \tag{11}$$

where, $L = \begin{bmatrix} K(A, C^T) & e_1 \end{bmatrix}$, $N = \begin{bmatrix} K(B, C^T) & e_2 \end{bmatrix}$, let $z_2 = \begin{bmatrix} w^{(2)} & b^{(2)} \end{bmatrix}^T$, then $z_2 = -(N^T N)^{-1} L^T \gamma$. The above analysis shows that once z_1 and z_2 are determined, two categorical hyperplanes are also determined.

3.3. Wavelet Kernel and Wavelet Twin Support Vector Machine

In the training process of nonlinear support vector machine, the selection of kernel function plays a crucial role in the classification effect of learning methods, so the selection of kernel function often affects the success or failure of the whole model application. In order to further improve the generalization ability of CNN-TWSVM model, the wavelet kernel function is introduced to build the wavelet twin support vector machine, so that the model can be better applied to small sample classification.

Kernel functions in support vector machines need to satisfy Mercer's theorem. Literature [41] proposed the translational invariance of kernel function based on Mercer's condition: $K(x, x') = K(x - x')$. However, it is difficult to decompose the translational invariant kernel into a product of two functions and prove that they are support vector kernels. For this reason, Zhang *et al.* [42] proposed a necessary and sufficient condition for translational invariant kernel.

Lemma 1: A translation invariant kernel $K(x, x') = K(x - x')$ is an admissible SV kernels if and only if the Fourier transform

$$F[K](\omega) = (2\pi)^{-N/2} \int_{R^N} \exp(-j(\omega \cdot x)) K(x) dx \tag{12}$$

is non-negative.

Lemma 2: Let $\varphi(x)$ be a mother wavelet, and let a and b denote the dilation and translation, respectively. $x, a, b \in R$. If $x, x' \in R^N$, then dot-product wave-

let kernels are

$$K(x, x') = \prod_{i=1}^N \varphi\left(\frac{x_i - b_i}{a}\right) \varphi\left(\frac{x'_i - b'_i}{a}\right) \quad (13)$$

and translation-invariant wavelet kernels that satisfy the translation invariant kernel theorem are

$$K(x, x') = \prod_{i=1}^N \varphi\left(\frac{x_i - x'_i}{a}\right). \quad (14)$$

In this paper, we select the Mexican Hat wavelet function as the translation-invariant wavelet kernel function. It is

$$\varphi(x) = (1 - x^2) \exp\left(-\frac{1}{2}x^2\right). \quad (15)$$

Lemma 3: The Mexican Hat wavelet kernel function that satisfies the translation-invariant kernel conditions is

$$K(x, x') = \prod_{i=1}^M \left(1 - \left(\frac{x_i - x'_i}{a_i}\right)^2\right) \exp\left(-\frac{1}{2} \left(\frac{x_i - x'_i}{a_i}\right)^2\right) \quad (16)$$

Theorem 1: The following formula is also a wavelet kernel function that satisfies the translation-invariant kernel conditions:

$$K(x, x') = \left(M - \sum_{i=1}^M \left(\frac{x_i - x'_i}{a_i}\right)^2\right) \exp\left(-\frac{1}{2} \sum_{i=1}^M \left(\frac{x_i - x'_i}{a_i}\right)^2\right) \quad (17)$$

According to *Theorem 1*, it can be proved that the above formula [43] satisfies Mercer's Theorem if and only if its Fourier transform is non-negative, forming a support vector kernel.

3.4. Traffic Sign Recognition Model Based on CNN-TWSVM

Based on the feature representation ability of CNN network and the advantages of support vector machine in dealing with small sample classification, a CNN-SVM model framework was established. In this paper, the parameters of each layer are retained by the trained CNN, and the convolutional layer and the pooling layer are intercepted to build a feature extraction model. The research shows that when the training data of the problem domain is not enough to achieve the optimal training effect of the CNN network, the feature expression gain can be obtained in a more perfect data set by means of pre-training, so as to ensure that the inner layer of the newly established CNN model still has a strong ability of feature extraction and representation. After normalization processing, the output of feature extraction model is used as the input of SVM classifier to complete the classification of local tasks.

The training process of CNN convolution is the process of increasing the degree of linear separability of a linearly non-fractionable data. After the characteristics extracted by the trained CNN network, the samples gradually tend to be linearly separable. Aiming at this kind of problem, combined with the advantage

that SVM only uses support vector samples to classify, the classification accuracy and generalization performance can be further improved. Therefore, the performance of SVM classifier will greatly affect the generalization ability and final recognition rate of the model. In this paper, twin support vector machines with higher classification efficiency are used to complete the classification. Aiming at the particularity of the image processing problem, the wavelet kernel function is introduced to deal with the nonlinear image classification problem, and the recognition rate of the model is further improved.

The feature extraction model in this paper is based on the VGG-16 network, and the network parameters are obtained by the classification training of the general image data set ImageNet. On the basis of preserving the internal network structure of the model, partial normalization processing is added to reduce over-fitting, and convolution layer and pooling layer are additionally added to ensure the feature representation ability while improving the generalization ability of the original model, which is more suitable for the classification of local tasks. The specific structure is shown in the figure below. Finally, the features extracted from the model will be classified by TWSVM. The specific construction steps of CNN-TWSVM model are as follows:

Step 1: The pre-trained VGG16 network was taken as the basic model, and the last three full connection layers and Softmax layers of the original model were removed.

Step 2: After the network obtained in Step 1, the full connection layer specially used for traffic sign classification is reconstructed, and the Softmax layer, which is the same as the traffic sign category, is set as the classifier.

Step 3: Fix the network parameters in Step 1, train the network constructed in Step 2 on the GTSRB data set, and save the network.

Step 4: The network trained in step 3 is used to extract the output of the next layer of Softmax layer and input it into the TWSVM classifier as input features for training and testing.

The schematic diagram of the traffic sign method proposed in this paper is shown in **Figure 2**, which mainly includes two processes: the establishment of feature mapping model by using CNN transfer learning and feature classification based on TWSVM.

Compared with the standard VGG16 model, three full-connection layers and the Softmax layer were deleted, and the full-connection layer structure for traffic sign recognition was reconstructed. As shown in the figure, four full-connection layers were added, with the number of neurons being 512, 256, 128 and 64 respectively. A dropout layer was added after each full-connection layer, with loss parameter Settings of 0.5, 0.6, 0.55 and 0.5 respectively. By fixing the network parameters of the convolutional layer, the whole network is retrained, and the feature mapping model is established. Secondly, the TWSVM model is used to replace the Softmax layer, and the features extracted through the network are input into the TWSVM to achieve the classification of traffic signs.

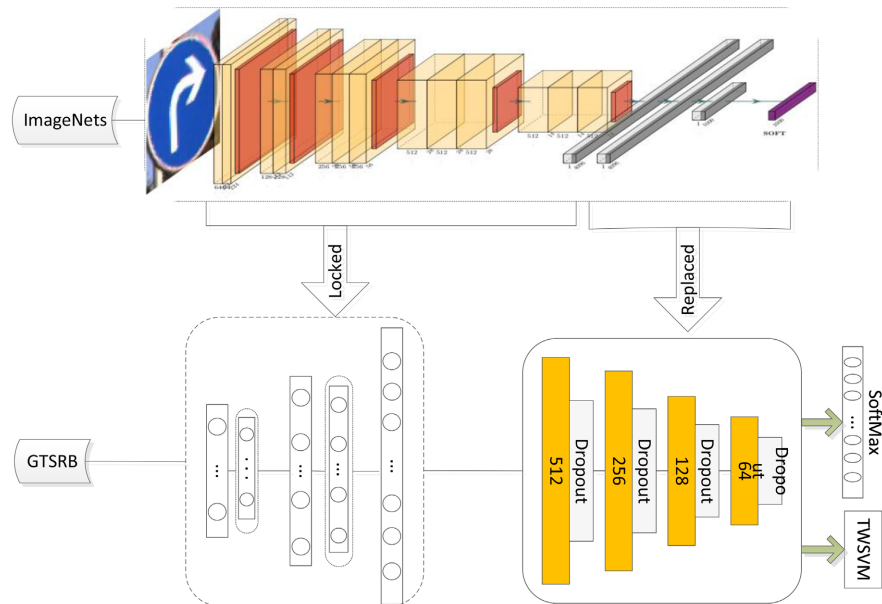


Figure 2. Schematic diagram of CNN-TWSVM model.

4. Experimental Results

4.1. Image Database and Data Processing

The CNN-TWSVM hybrid model based on wavelet kernel function proposed in this paper is trained on the classical German Traffic Sign Data Set (GTSRB). Then the GTSRB data set is used to build a test data set to test the completed training model. GTSRB contains nearly 39,000 pictures of 43 types of traffic signs. The pictures of each type of traffic signs contain traffic scenes in various complex situations, as shown in **Figure 3**, which is suitable for training the network in this paper. During the construction of the test data set, in order to ensure the structural balance of the data, 100 pictures of each category were randomly selected from the GTSRB data set, and a total of 4300 pictures constituted the test set.

4.2. Test Results and Analysis

In order to verify the practical application effect of the model proposed in this paper, this section conducts comparative experiments on two traffic sign data sets respectively. The mixed model constructed based on transfer learning idea in this paper is trained on GTSRB data set, and the training effect and confusion matrix of CNN-SVM and CNN-TWSVM models are compared under Gaussian kernel function and wavelet kernel function. Then the model is tested on the test data set to evaluate the application effect of the new model. For the super parameter optimization problem in the support vector machine model, the optimal parameters are sought by using the optimization algorithm of grid search, and then the optimal parameters are substituted back to the model for the second training. In the training process of CNN-TWSVM model, in addition to replacing the last classification layer of the model with the new support vector classifier

model, the remaining network parameters will apply the training results of CNN-Softmax model on the GTSRB data set. The training behavior for the new model only trains the parameters inside the classification model at the last level, which is equivalent to applying the idea of transfer learning for the second time. Using this training method, we can fully compare the effects of different classification models.

Figure 4 shows the training results of the convolutional neural network for traffic sign classification constructed in this paper on the GTSRB dataset respectively. On the GTSRB data set, the training accuracy of CNN-Softmax is stable at 96.80%, the verification accuracy is stable at 96.70%, the training loss is stable at about 0.125, and the verification loss is stable at about 0.23. This indicates that



Figure 3. Sample instance of GTSRB dataset.

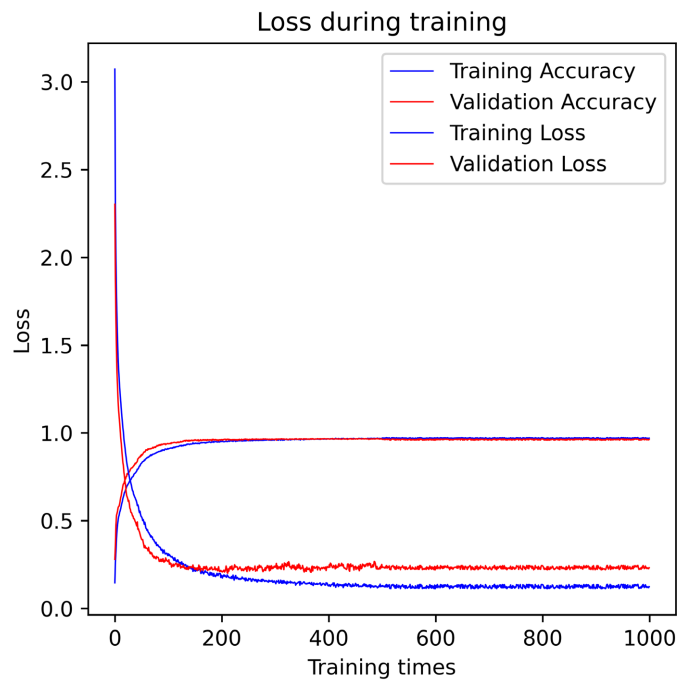


Figure 4. Training and validation accuracy and loss.

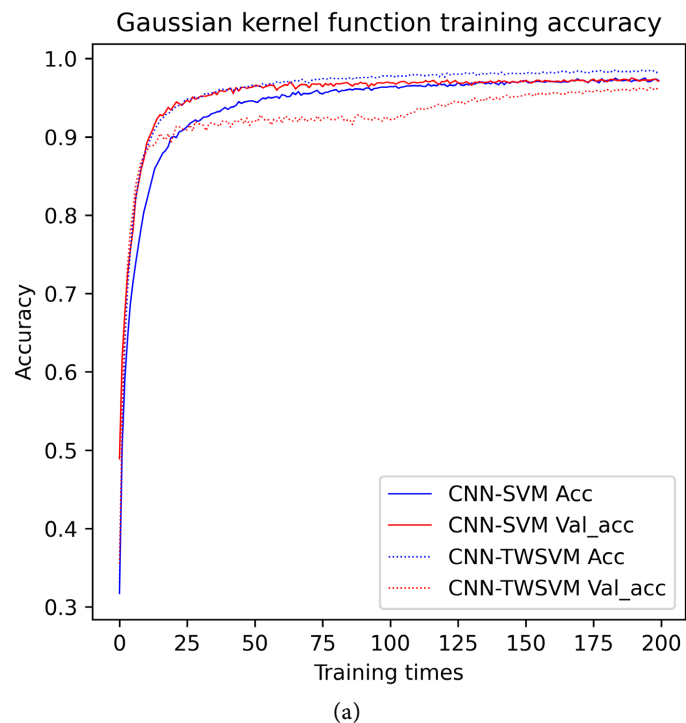
the training of the network is good, and the training loss can be further reduced by adjusting the network structure to make the network fit better. In the test data, the test accuracy reached 98.30%.

In order to test the performance of each model when using different kernel functions on this data set, the traditional CNN-SVM model and Gaussian kernel function were introduced to carry out comparative tests. The training accuracy and verification accuracy of CNN-SVM model and CNN-TWSVM under Gaussian kernel function and wavelet kernel function are shown in **Figure 5(a)** and **Figure 5(b)**.

Table 1 shows the stable training accuracy and verification accuracy of the two comparison models under different kernel functions, as well as the accuracy displayed on the test set. It can be seen from the data in the table that the new hybrid model proposed in this paper can effectively reduce over-fitting, and the generalization ability of the model can be further improved with the cooperation of the wavelet kernel function.

Table 1. Model comparison.

| | | CNN-SVM | CNN-TWSVM |
|---------|---------------------|---------|-----------|
| RBF | accuracy | 97.23% | 98.12% |
| | Validation accuracy | 97.16% | 98.06% |
| | Test accuracy | 96.70% | 97.32% |
| WAVELET | accuracy | 97.42% | 98.74% |
| | Validation accuracy | 97.25% | 98.71% |
| | Test accuracy | 96.77% | 99.12% |



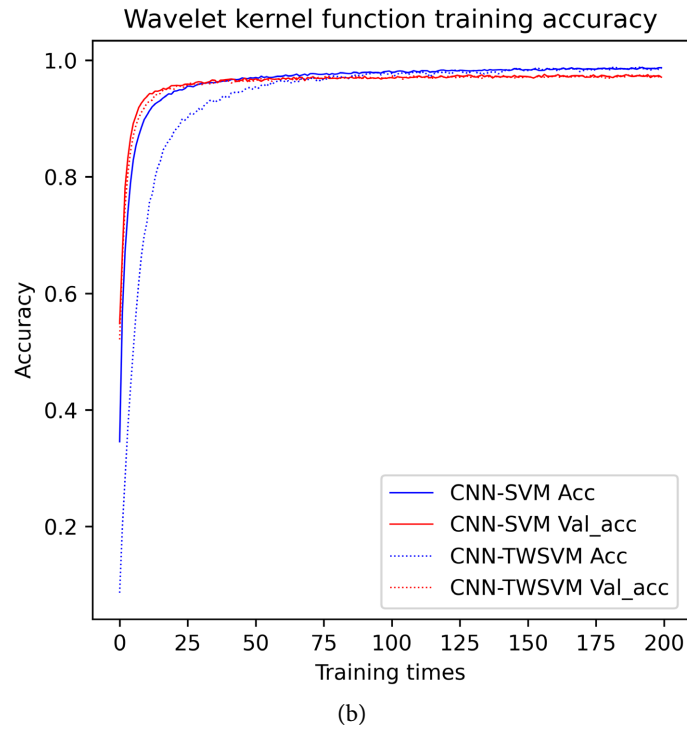


Figure 5. Training and validation accuracy of CNN-SVM and CNNTWSVM. (a) Gaussian kernel. (b) Wavelet kernel.

4.3. Accuracy Evaluation Index

Mean accuracy (MAP) [44] is a commonly used index to evaluate accuracy in the field of target detection and recognition. First, calculate precision and recall:

$$\text{precision} = \frac{TP}{TP + FP}, \quad \text{recall} = \frac{TP}{TP + FN}$$

where, TP , FP and FN respectively represent correctly identified positive samples, wrongly identified positive samples and wrongly identified negative samples. AP is defined as follows:

$$AP_i = \sum_{k=1}^N p(k) \Delta r(k) \tag{18}$$

where, $p(k)$ represents the accuracy rate corresponding to point k of recall rate change; $\Delta r(k)$ represents the change amount of recall rate corresponding to change point k ; N represents the number of recall rate change points; Different categories have different values of AP , and i is the index value of the category. $p(k)$ is defined as follows:

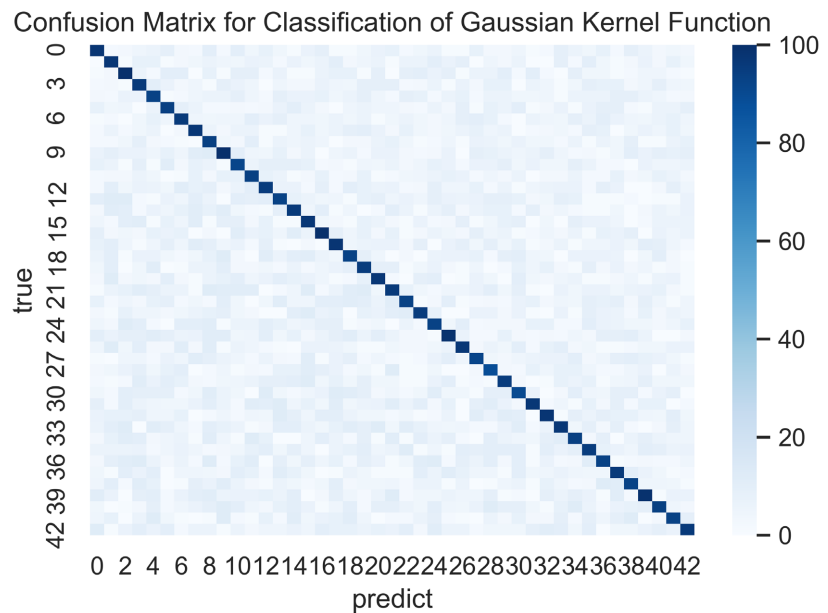
$$p(k) = \max_{k:\tilde{k} \geq k} p(\tilde{k}) \tag{19}$$

Here, the accuracy $p(k)$ corresponding to recall rate k is the maximum accuracy of any recall rate in $\tilde{k} \geq k$, and $p(\tilde{k})$ is the accuracy corresponding to recall rate \tilde{k} . mAP is the mean value of AP values of all categories, defined as follows:

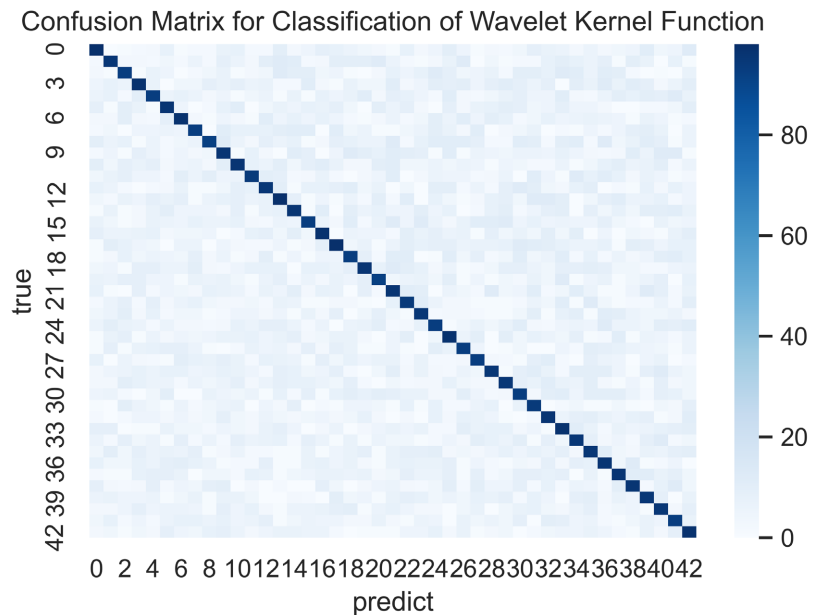
$$mAP = \frac{1}{m} \sum_{i=1}^m AP_i \quad (20)$$

where m is the number of categories.

Figure 6 shows the confusion matrix of CNN-TWSVM model under Gaussian kernel function and wavelet kernel function on the test set. **Table 2** shows the mean accuracy of each kernel function on the test set for different SVM models.



(a)



(b)

Figure 6. Comparison of confusion matrix. (a) RBF confusion matrix; (b) WAVELET confusion matrix.

Table 2. MComparison of average accuracy.

| Model | <i>mAP</i> |
|---------------------|------------|
| CNN-softmax | 87.23% |
| CNN-SVM (rbf) | 90.37% |
| CNN-SVM (wavelet) | 91.24% |
| CNN-TWSVM (rbf) | 92.42% |
| CNN-TWSVM (wavelet) | 93.54% |

Experimental results show that the proposed model in the CNN-TWSVM model with Gaussian kernel function and wavelet kernel function is better than the traditional CNN-SVM model and CNN-Softmax model in both the test accuracy and average precision mean value. This indicates that the improved model can effectively avoid over-fitting in the task of traffic sign recognition and can improve the generalization ability of the model.

5. Concluding Remarks

Based on the idea of migration learning, this paper reconstructs the feature extraction model for traffic sign recognition based on the VGG16 network. The traditional CNN-SVM model is improved, and the twin support vector machine with higher computational efficiency and classification efficiency is introduced into the model, and the wavelet kernel function is used to replace the Gaussian kernel function widely used before to process data with image characteristics. The experiment verified the rationality of the hybrid model based on the CNN migration learning feature mapping model and the TWSVM classification model fusion, and the wavelet kernel function used can effectively improve the generalization ability of the model, making it have excellent recognition in small sample recognition tasks Accuracy.

Based on the improvement ideas of the convolutional neural network model in this article, it is not difficult to notice that in the process of inputting convolutional features into the classifier, there is also a process of tiling three-dimensional tensor data into vector features. This undoubtedly destroys the structure of the data and loses the structural information between the original three-dimensional data. In the face of this kind of damage, although a certain degree of correction can be obtained by improving the subsequent classification model, the damage to the structural information is irreversible. Therefore, in future research, it will be possible to focus on the idea of maintaining the structural information of the convolution feature. For example, the introduction of support tensor machine to retain more structural information. Among them, how to maintain more tensor structure information itself is also a hot research topic. And this will also become our research work for some time in the future.

Acknowledgements

In the process of completing the study, the author thanks for being supported by

two fund projects. They are China National Tobacco Corporation Yunnan Science and technology planning project “Application Research of data resource mining based on multi-objective decision” (No. 80026091555) and Yunnan University of Finance and economics graduate innovation project fund projects “Research on medical image classification of lung cancer based on the new CNN-SVM hybrid model” (grant No. 2021YUFEYC074) respectively.

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] De La Escalera, A., Moreno, L.E., Salichs, M.A. and Armingol, J.M. (1997) Road Traffic Sign Detection and Classification. *IEEE Transactions on Industrial Electronics*, **44**, 848-859. <https://doi.org/10.1109/41.649946>
- [2] Zhu, S., Liu, L. and Lu, X. (2006) Color-Geometric Model for Road Traffic Sign Recognition. *The Proceedings of the Multiconference on “Computational Engineering in Systems Applications”*, Beijing, 4-6 October 2006, 2028-2032. <https://doi.org/10.1109/CESA.2006.4281972>
- [3] Bahlmann, C., Zhu, Y., Ramesh, V., Pellkofer, M. and Koehler, T. (2005) A System for Traffic Sign Detection, Tracking, and Recognition Using Color, Shape, and Motion Information. *IEEE Proceedings of Intelligent Vehicles Symposium*, Las Vegas, 6-8 June 2005, 255-260. <https://doi.org/10.1109/IVS.2005.1505111>
- [4] Miao, X., Li, S.M., Shen, Y. and Yu, H.B. (2013) Real Time Recognition Method of Traffic Signs in Complex Environment. *Journal of Jiangsu University: Natural Science Edition*, **34**, 514-518.
- [5] Xu, S.Q. (2009) Outdoor Traffic Sign Detection and Shape Recognition. *Chinese Journal of Image and Graphics*, **14**, 707-711.
- [6] Gil-Jimenez, P., Lafuente-Arroyo, S., Gomez-Moreno, H., Lopez-Ferreras, F. and Maldonado-Bascon, S. (2005) Traffic Sign Shape Classification Evaluation. Part II.FFT Applied to the Signature of Blobs. *IEEE Proceedings of Intelligent Vehicles Symposium*, Las Vegas, 6-8 June 2005, 607-612. <https://doi.org/10.1109/IVS.2005.1505170>
- [7] Gao, X.W., Podladchikova, L., Shaposhnikov, D., Hong, K. and Shevtsova, N. (2006) Recognition of Traffic Signs Based on Their Colour and Shape Features Extracted Using Human Vision Models. *Journal of Visual Communication and Image Representation*, **17**, 675-685. <https://doi.org/10.1016/j.jvcir.2005.10.003>
- [8] Broggi, A., Cerri, P., Medici, P., Porta, P.P. and Ghisio, G. (2007) Real Time Road Signs Recognition. 2007 *IEEE Intelligent Vehicles Symposium*, Istanbul, 13-15 June 2007, 981-986. <https://doi.org/10.1109/IVS.2007.4290244>
- [9] Sun, G.M., Wang, J., Yu, G.Y., Li, G. and Xu, L. (2010) Detection and Recognition of Traffic Signs in Natural Background. *Journal of Beijing University of Technology*, No. 10, 1337-1343.
- [10] Moutarde, F., Bargeton, A., Herbin, A. and Chanussot, L. (2007) Robust On-Vehicle Real-Time Visual Detection of American and European Speed Limit Signs, with a Modular Traffic Signs Recognition System. 2007 *IEEE Intelligent Vehicles Symposium*, Istanbul, 13-15 June 2007, 1122-1126.

- <https://doi.org/10.1109/IVS.2007.4290268>
- [11] Shams, M.M., Kaveh, H. and Safabakhsh, R. (2015) Traffic Sign Recognition Using an Extended Bag-of-Features Model with Spatial Histogram. 2015 *Signal Processing and Intelligent Systems Conference (SPIS)*, Tehran, 16-17 December 2015, 189-193. <https://doi.org/10.1109/SPIS.2015.7422338>
- [12] Dalal, N. and Triggs, B. (2005) Histograms of Oriented Gradients for Human Detection. 2005 *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, Vol. 1, San Diego, 20-25 June 2005, 886-893. <https://doi.org/10.1109/CVPR.2005.177>
- [13] Lowe, D.G. (1999) Object Recognition from Local Scale-Invariant Features. *Proceedings of the 7th IEEE International Conference on Computer Vision*, Vol. 2, Kerkyra, 20-27 September 1999, 1150-1157. <https://doi.org/10.1109/ICCV.1999.790410>
- [14] Hu, X., Zhu, X., Li, D. and Li, H. (2010) Traffic Sign Recognition Using Scale Invariant Feature Transform and SVM. *A Special Joint Symposium of ISPRS Technical Commission IV and AutoCarto in Conjunction with ASPRS/CaGIS Fall Specialty Conference*, Orlando, 15-19 November 2010, 16-19.
- [15] Bay, H., Tuytelaars, T. and Van Gool, L. (2006) Surf: Speeded up Robust Features. *European Conference on Computer Vision*, Graz, 7-13 May 2006, 404-417. https://doi.org/10.1007/11744023_32
- [16] Hoferlin, B. and Zimmermann, K. (2009) Towards Reliable Traffic Sign Recognition. 2009 *IEEE Intelligent Vehicles Symposium*, Xi'an, 3-5 June 2009, 324-329. <https://doi.org/10.1109/IVS.2009.5164298>
- [17] Creusen, I.M., Wijnhoven, R.G., Herbschleb, E. and de With, P.H.N. (2010) Color Exploitation in Hog-Based Traffic Sign Detection. 2010 *IEEE International Conference on Image Processing*, Hong Kong (China), 26-29 September 2010, 2669-2672. <https://doi.org/10.1109/ICIP.2010.5651637>
- [18] Qingsong, X., Juan, S. and Tiantian, L. (2010) A Detection and Recognition Method for Prohibition Traffic Signs. 2010 *International Conference on Image Analysis and Signal Processing*, Zhejiang, 9-11 April 2010, 583-586. <https://doi.org/10.1109/IASP.2010.5476048>
- [19] Xie, Y., Liu, L.F., Li, C.H. and Qu, Y.Y. (2009) Unifying Visual Saliency with HOG Feature Learning for Traffic Sign Detection. 2009 *IEEE Intelligent Vehicles Symposium*, Xi'an, 3-5 June 2009, 24-29. <https://doi.org/10.1109/IVS.2009.5164247>
- [20] Huang, Z., Yu, Y., Gu, J. and Liu, H. (2016) An Efficient Method for Traffic Sign Recognition Based on Extreme Learning Machine. *IEEE Transactions on Cybernetics*, **47**, 920-933. <https://doi.org/10.1109/TCYB.2016.2533424>
- [21] Abedin, M.Z., Dhar, P. and Deb, K. (2016) Traffic Sign Recognition Using Hybrid Features Descriptor and Artificial Neural network Classifier. 2016 *19th International Conference on Computer and Information Technology (ICCIIT)*, Dhaka, 18-20 December 2016, 457-462. <https://doi.org/10.1109/ICCITECHN.2016.7860241>
- [22] Stallkamp, J., Schlipsing, M., Salmen, J. and Igel, C. (2012) Man vs. Computer: Benchmarking Machine Learning Algorithms for Traffic Sign Recognition. *Neural networks*, **32**, 323-332. <https://doi.org/10.1016/j.neunet.2012.02.016>
- [23] CireAn, D., Meier, U., Masci, J. and Schmidhuber, J. (2012) Multi-Column Deep Neural Network for Traffic Sign Classification. *Neural Networks*, **32**, 333-338. <https://doi.org/10.1016/j.neunet.2012.02.023>
- [24] Sermanet, P. and LeCun, Y. (2011) Traffic Sign Recognition with Multi-Scale Con-

- volutional Networks. *The 2011 International Joint Conference on Neural Networks*, San Jose, 31 July-5 August 2011, 2809-2813.
<https://doi.org/10.1109/IJCNN.2011.6033589>
- [25] Jin, J., Fu, K. and Zhang, C. (2014) Traffic Sign Recognition with Hinge Loss Trained Convolutional Neural Networks. *IEEE Transactions on Intelligent Transportation Systems*, **15**, 1991-2000. <https://doi.org/10.1109/TITS.2014.2308281>
- [26] Yin, S.H., Deng, J.C., Zhang, D.W. and Du, J.Y. (2017) Traffic Sign Recognition Based on Deep Convolutional Neural Network. *CCF Chinese Conference on Computer Vision*, Tianjin, 11-14 October 2017, 685-695.
https://doi.org/10.1007/978-981-10-7299-4_57
- [27] Hu, W.Z., Zhuo, Q., Zhang, C.S. and Li, J.K. (2017) Fast Branch Convolutional Neural Network for Traffic Sign Recognition. *IEEE Intelligent Transportation Systems Magazine*, **9**, 114-126. <https://doi.org/10.1109/MITS.2017.2709780>
- [28] Zeng, Y., Xu, X., Fang, Y. and Zhao, K. (2015) Traffic Sign Recognition Using Extreme Learning Classifier with Deep Convolutional Features. *The 2015 International Conference on Intelligence Science and Big Data Engineering (ISCIDE 2015)*, Vol. 9242, Suzhou, 14-16 June 2015, 272-280.
https://doi.org/10.1007/978-3-319-23989-7_28
- [29] Huang, G.B., Zhu, Q.Y. and Siew, C.K. (2004) Extreme Learning Machine: A New Learning Scheme of Feedforward Neural Networks. 2004 *IEEE International Joint Conference on Neural Networks (IEEE Cat. No. 04CH37541)*, Vol. 2, Budapest, 25-29 July 2004, 985-990. <https://doi.org/10.1109/IJCNN.2004.1380068>
- [30] Mathias, M., Timofte, R., Benenson, R. and Van Gool, L. (2013) Traffic Sign Recognition—How Far Are We From the Solution? *The 2013 international Joint Conference on Neural Networks (IJCNN)*, Dallas, 4-9 August 2013, 1-8.
<https://doi.org/10.1109/IJCNN.2013.6707049>
- [31] Yang, Y., Luo, H., Xu, H. and Wu, F. (2015) Towards Real-Time Traffic Sign Detection and Classification. *IEEE Transactions on Intelligent Transportation Systems*, **17**, 2022-2031. <https://doi.org/10.1109/TITS.2015.2482461>
- [32] Lian, Z., Jing, X., Sun, S. and Huang, H. (2016) Frequency Selective Convolutional Neural Networks for Traffic Sign Recognition. 2016 *IEEE 83rd Vehicular Technology Conference (VTC Spring)*, Nanjing, 15-18 May 2016, 1-5.
<https://doi.org/10.1109/VTCSpring.2016.7504251>
- [33] Jang, C., Kim, H., Park, E. and Kim, H. (2016) Data Debaised Traffic Sign Recognition Using MSERs and CNN. 2016 *International Conference on Electronics, Information, and Communications (ICEIC)*, Danang, 27-30 January 2016, 1-4.
<https://doi.org/10.1109/ELINFOCOM.2016.7562938>
- [34] Deng, Z.D. and Zhou, L.P. (2017) Detection and Recognition of Traffic Planar Objects Using Colorized Laser Scan and Perspective Distortion Rectification. *IEEE Transactions on Intelligent Transportation Systems*, **19**, 1485-1495.
<https://doi.org/10.1109/TITS.2017.2723902>
- [35] Gudigar, A., Chokkadi, S., Raghavendra, U. and Acharya, U.R. (2017) Multiple Thresholding and Subspace Based Approach for Detection and Recognition of Traffic Sign. *Multimedia Tools and Applications*, **76**, 6973-6991.
<https://doi.org/10.1007/s11042-016-3321-6>
- [36] Tan, M., Wang, B., Wu, Z., Wang, J. and Pan, G. (2016) Weakly Supervised Metric Learning for Traffic Sign Recognition in a LIDAR-Equipped Vehicle. *IEEE Transactions on Intelligent Transportation Systems*, **17**, 1415-1427.
<https://doi.org/10.1109/TITS.2015.2506182>

- [37] Chen, Z.S. and Zhang, D.F. (2017) Road Marking Recognition Based on Gaussian-Hermite Moment. *Proceedings of the International Conference on Advances in Image Processing*, Bangkok, 25-27 August 2017, 24-27.
<https://doi.org/10.1145/3133264.3133279>
- [38] Lu, X., Wang, Y., Zhou, X., Zhang, Z. and Ling, Z. (2017) Traffic Sign Recognition via Multi-Modal Tree-Structure Embedded Multi-Task Learning. *IEEE Transactions on Intelligent Transportation Systems*, **18**, 960-972.
<https://doi.org/10.1109/TITS.2016.2598356>
- [39] Arcos-Garcia, A., Soilan, M., Alvarez-Garcia, J.A. and Riveiro, B. (2017) Exploiting Synergies of Mobile Mapping Sensors and Deep Learning for Traffic Sign Recognition Systems. *Expert Systems with Applications*, **89**, 286-295.
<https://doi.org/10.1016/j.eswa.2017.07.042>
- [40] Zhang, J., Wang, W., Lu, C., Wang, J. and Sangaiah, A.K. (2020) Lightweight Deep Network for Traffic Sign Classification. *Annals of Telecommunications*, **75**, 369-379. <https://doi.org/10.1007/s12243-019-00731-9>
- [41] Smola, A.J., Schölkopf, B. and Müller, K.R. (1998) The Connection between Regularization Operators and Support Vector Kernels. *Neural Networks*, **11**, 637-649.
[https://doi.org/10.1016/S0893-6080\(98\)00032-X](https://doi.org/10.1016/S0893-6080(98)00032-X)
- [42] Zhang, L., Zhou, W.D. and Jiao, L.C. (2004) Wavelet Support Vector Machine. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, **34**, 34-39.
<https://doi.org/10.1109/TSMCB.2003.811113>
- [43] Ding, S., Wu, F. and Shi, Z. (2014) Wavelet Twin Support Vector Machine. *Neural Computing and Applications*, **25**, 1241-1247.
<https://doi.org/10.1007/s00521-014-1596-y>
- [44] Aurand, A.M., Dufour, J.S. and Marras, W.S. (2017) Accuracy Map of an Optical Motion Capture System with 42 or 21 Cameras in a Large Measurement Volume. *Journal of Biomechanics*, **58**, 237-240.
<https://doi.org/10.1016/j.jbiomech.2017.05.006>