

TRAFFIC SIGN DETECTION USING CNN

Chatrathy Pavan Tirumala Sai Gopal

Reg No: 11705122

KM041

INT 248

Roll No: 09

Abstract:

The most commonly used convolutionary neural networks are Deep learning algorithms for classification of traffic signals until date[1] But they struggle to capture the images' pose, vision, orientation, This paper is due to the intrinsic failure of the max pooling sheet. Proposes a new approach using deep detection of traffic signals Capsule networks, a learning architecture that achieves Excellent results on the German sign for traffic Dataset, dataset. The network of capsules consists of capsules that are a group of Neurons that describe an object's instantiating parameters, such as By using the dynamic routing and route, the pose and orientation[2] Algorithms by agreement. Unlike previous approaches Extraction of manual functions, several deep neural networks with Our system removes manualwork and several criteria, The spatial variances provide resistance. It is possible to trick CNNs' Different adversary attacks[3] and capsule networks can easily be used. Overcoming those attacks by intruders and providing more Traffic sign detection reliability for autonomous sign detection. Capsule network has accomplished 94 percent precision on German Traffic Sign Identification Dataset of Benchmarks (GTSRB)

Keywords : CNN, Capsule Net, Pose, Traffic sign, Dataset

I . INTRODUCTION

Detection of traffic signs is a task in the real world that requires a lot of Constraints and complications. Even a small volume Misclassification of the traffic sign will result in catastrophe. And it may also contribute to the loss of life. It is Implemented in different sophisticated schemes of driver assistance A camera is mounted on the camera and in autonomous vehicles. The car dashboard records the video in real time and Feed that is sampled in frames and fed to profound learning model which is deployed within an automotive embedded Board. As the vehicle is driven in different ways, Ecosystems, lighting conditions, speeds and geographies are The deep learning algorithm needs to be stable and stable. Reliable at all times. The camera will catch the sign

in traffic Various orientations and poses, but the algorithm ought to be The right sign[4] and capsule networks are capable of recognising To tackle this, the optimal deep learning algorithm Problem.

Convolutionary neural networks are commonly used for all State of the art neural network algorithms for deep learning[5] In most of the tasks linked to the picture, Convolution captures the Spatial image knowledge using the role of the kernel in Layer convolution. ACNN is made up of input, output and concealed Of layers. Furthermore, the secret layers consist of convolutionary ones, Pooling, completely connected and layers of normalisation. CNNs for They do very well in image-related tasks, but they have Few simple constraints and draw backs. CNN fails to[6] Capture the relative orientation and spatial relationships. CNN can get confused easily by image orientation or by change in pose.

Details about the pose can be rotation, thickness, skew, precision, Place object. The biggest downside to the Max pooling is CNN because the spatial hierarchies between them can not be propagated Simple and complex objects which contribute to invariance and invariance It prevents them from capturing the pose and spatial The relationships in the picture between the pixels. CNN uses the The max pooling layer that samples the data and decreases the data Spatial data information which is passed on to the next layer

This downside capsnet architecture has been invented to solve The state of the art success on MNIST reached Digits dataset[2] and better results have been obtained than CNNs on Dataset Multi MNIST.

II . RELATED WORK

It's hard to compare the research work performed in the past. The area of traffic sign detection due to the comprehensive Study efforts carried out by scientists in this field and the use of Various forms of datasets to solve various problems Detection, classification and monitoring are included. Tasks linked.

A. Using computer vision feature extraction methods

This is one of the early methods in which a great deal of Computer algorithms and methodologies have been proposed by Before the invention of robots, vision scientists Learning. Techniques such as Gradients with Histogram Orientation Initially, HOG)[7] was popularised for the identification of Pedestrian. In this process, colour image gradients are Calculated together with various normalised, weighted Histograms. Scale

Invariant Transform of Features (SIFT)[8] The approach for classification and sliding window was used The system was used to carry out both grouping and grouping. Simultaneous detection assignments.

B. Using machine learning

Many types of machine learning algorithms[9], such as Help vector machines, study of linear discriminants[10], The classifiers of the ensemble, random woods, and kd-trees[11] were Used for the designation of traffic signals. The basis of Linear Discriminate Analysis(LDA)[5] is Maximum estimate of class membership posteriori. Class membership Densities are known to have Gaussian and Gaussian multi-variate densities Popular matrix co-variance. Random Forest is a tool for ensemble classification[1], which is Based on the set of random decision trees that are not pruned. Each one The decision tree is constructed using training data taken randomly. For All decision trees compare classification test data from all decision trees.And the performance of the classification is based on majority voting, Consideration of the decisions of all decision trees of the majority.

Help Vector Machines(SVM) is a category The algorithm that classifies the knowledge by dividing the nA hyper plane with a dimensional data plane for Classifying[9].

SVM is also able to distinguish non-linearly.By converting the classification plane to the classification plane, scattered data Higher dimensions using a non-linear function of the kernel using For its implementation, a mechanism called kernel trick.

Machine learning strategies[12] were unable to deal with Variable aspect ratio and pictures with variable size and characteristics It must be hand engineered manually, which is very time consuming. A method vulnerable to consumption and error.

C. Using deep learning

To resolve the drawbacks mentioned above, Fresh implementations of conventional methodologies based on The previous techniques[13] have been replaced by deep learning algorithms. In recent years, there has been an increase in computing capacity and Availability of structured data sets and access to enormous Data quantity. The condition of convolutionary neural networks is Best accuracy rate of the art algorithms. LENET The first CNN architecture for traffic sign architecture[14] was Classifying. Biologically motivated, convolutionary neural networks are Architecture of the multi-stage neural network that learns the Automatically Invariant Functions. The philtre for each stage is Bank(convolution) layer, layer of non-linear transformation, spatial transformation, The layer of pooling[15]. The layer of spatial pooling deceases the Spatial details and functions in visual form as a complex cell Cortex. For preparation, a gradient descent based optimizer is used In order to minimise the loss function, update each philtre. The output of all the layers is fed to the classifier for improving the accuracy of classification.

III . DATASET

The German Benchmark for Traffic Sign Identification (GTSRB) It is a publicly accessible dataset that has been identified and visualised.Dataset which is built from a 10-hour driving video on The video is shot using Prosilica on various roads in Germany.Camera GC 1380CH with 25 fps framerate and traffic The NISYS Advanced extraction of signs is done using the System for Growth and Analysis(ADAF)[16] module Software framework that is based. The redundant and repetitive after washing and removing the The dataset frames are reduced to a total of 51,840 of the 433 images. Oh, lessons. In the dataset, all the images have a scale of 32 * 32 and theThe overall dataset is split into data for training and data for research.There are 39,209 images and 12,630 images as training info. As data for assessments.



Fig. 1. Sample images from GTSRB dataset

IV. CAPSULE NET ARCHITECTURE

Capsule networks are made of capsules instead of neurons.Capsule[17] is a community of neural artificial networks performing On their inputs and complicated internal computations Encapsulate a tiny vector with the data. Every single capsule Captures the object's relative location and whether the object is The location is altered and the direction of the output vector is changed. Changed[18], thus making them equivariant. Caps Net is made up of several layers and the first layer is Primary capsules, where each capsule receives a tiny capsule As data, part of the receptive field and attempts to detect the pose Unique pattern. The capsule's production is a vector and To ensure the performance, dynamic routing technique was used to ensure It is sent in the layer to the appropriate parent, who can be Subtracted from Fig. 2.

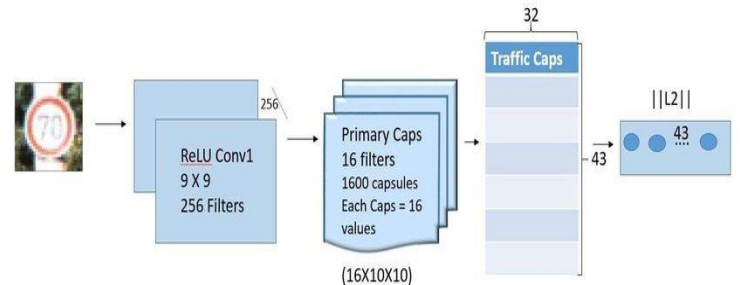


Fig. 2. Capsule sign architecture for traffic sign detection

A. Computation of capsule vector inputs and outputs

By multiplying the weight matrix(W_{ij}) with its own output vector(u_i), the capsule computes a prediction vector. For that specific capsule output, the coupling coefficient of that capsule's corresponding output increases the scalar product and prediction.

$$u_j|i = W_{ij}u_i$$

where $u_j|i$ =prediction vector, W_{ij} =weight matrix and u_i =output vector.

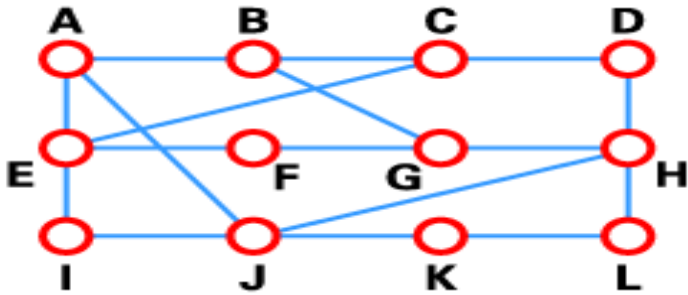
B. Squash Function

A non-linear activation mechanism called the squashing function is used in capsule networks. This function transforms the output vector's length into the capsule's likelihood of connecting to that object. It shrinks the long output vectors just below the length of one and the short output vectors near zero. Where s_j = Complete Input, v_j = Vector, $Output_j$ = Capsule outputj.

$$v_j = \frac{||s_j||^2}{1 + ||s_j||^2} \frac{s_j}{||s_j||}$$

C. Routing Algorithm

The CNN scalar output function detectors have been replaced by vector output capsules and max route pooling by arrangement. With the increase in hierarchy, the dimensionality of the capsule increases due to the change from position coding (encoded in continuous space) to rate coding (encoded in) and the high-level capsules represent entities that are complex and have more degree of freedom. This approach is successful by agreement than the maximum pooling used in the CNNs. Where the iterative dynamic routing defines S_j = summation matrix, u_j = prediction vector, and c_{ij} =coupling coefficients.



V. EXPERIMENTS

A. Data Preparation

The data is pickled from the disc. Pickling is the method of saving a file before writing it to the disc in a serialised format. The size of each image is 32 * 32 and there are a total of 34799 images in the training dataset and 12630 in the testing dataset.

B. preprocessing

With a random uniform distribution of 0.6 to 1.5, the image brightness is increased and image contrast is also enhanced with a random uniform distribution of 0.6 to 1.5. By replicating the available with data rotation ± 20 , shear range of 0.2, width shift range of 0.2, horizontal flip, which increased the size of the training dataset, leading to better output and regularising, thereby avoiding the issue of overfitting, the training dataset is increased five-fold. The dataset size will be increased to 34799 X 5 = 173995 images by increasing the current training dataset.

C. Network architecture

As part of primary capsules, the architecture used for traffic sign detection consists of the input layer and initially convolutional layers and the output vector of the primary capsule is sent to capsules of traffic signs.



Fig. 3. Visualizing the data

1) Input Layer: The input layer consists of images of input training and the dimension is equal to the images of total training.

2) Primary Capsule Layer: The primary capsule layer is the first layer that follows the input layer, and the first two convolution layers are used to measure the output. The first convolution layer consists of philtres of kernel size 9 and 256 and padding was not used. As a non-linear activation function, the rectified linear unit (ReLU) was used and a drop out of 0.7 was fixed to be optimal after checking with various values. Output is reshaped to obtain the output vectors of primary capsules. Because the primary capsule layer is entirely connected to the capsule layer of the traffic sign, the output vectors must be squashed using the squashing function. Small epsilon value is applied to the squash function to prevent the issue of the vanishing gradient during training.

IV . RESULT AND DISCUSSION

3) Traffic Sign Capsule Layer: To calculate the output of capsules of traffic signs, calculate the estimated output vectors for each primary pair of capsules of traffic signs and implement the direction by agreement algorithm. The capsule layer of the traffic sign consists of 43 capsules, each representing a special class of the German traffic sign dataset, each of which has a size of 32. The corresponding weights and output vectors of each capsule j in the second layer are predicted for each capsule I in the first layer.

D. Reconstruction

A decoder network is introduced to the capsule traffic sign network, which consists of a completely connected network layer that helps to reconstruct the input images by tuning the performance of the capsule traffic sign network. This feedback mechanism would allow the network to retain the information needed for the entire network to reconstruct the traffic sign.

1) Mask : Only the particular output vector corresponding to the expected traffic sign is sent for the reconstruction of the input traffic sign and all remaining outputs should be masked. The masking function is used during the training process to prevent all other output vectors. Using the one-hot function, the reconstruction mask is realised. Its value will be one for the target class and its value will be zero for all the other classes.

2)Decoder: The Decoder consists of a ReLU non-linear activation layer followed by a layer of sigmoid activation.

E. Losses

1) Margin Loss: The length of the instantiation output vector reflects whether or not the likelihood of the entity of the respective capsule occurs. Only if the traffic sign is present in the input picture does the digit class k have the longest vector output. The margin loss is separate for each traffic signal capsule k and it is given as

$$L_k = T_k \max(0, m^+ - \|v_k\|) + \lambda(1 - T_k) \max(0, \|v_k\| - m^-)$$

If a traffic sign of class k is present, the value of T_k is 1 and here $m^+ = 0.9$ and $m^- = 0.1$. λ is a regularisation parameter that prevents the learning of all traffic sign capsules from shrinking the operation vector.

2) Reconstruction Loss: It is the difference between the input image squares and the reconstructed image $R = (\text{Input image})^2 - (\text{Reconstructed image})^2$ where $R = \text{Reconstruction Loss}$ 3) Final Loss: The Final Loss is the sum scaled to a factor λ of margin loss and reconstruction loss that functions as a scaling factor and should be much less than one.

$$F = (\text{Margin Loss}) - \lambda(\text{Reconstruction Loss})$$

In contrast, where $F = \text{Final Loss}$, $\lambda = 0.0005$ Margin loss should always dominate the loss of reconstruction. If reconstruction loss is more in the final loss, then the model attempts to compare precisely the output image with the training dataset input image that contributes to the model overfitting to the training data.

Using the test data set of 12,630 test images, the model is tested. The accuracy is determined by the total number of traffic signs as the ratio of the correctly defined traffic signs.

The research dataset achieved an accuracy of 94 percent with a batch size of 50 and a final loss of 0.0311028 measured. The performance measurement is based on the right rate of classification (CCR) and binary loss (0 or 1), which ensures that the number of misclassifications is counted.

In terms of the number of samples for a particular class, the test set is unbalanced, but all classes are equally essential and have equal weight. In order to implement this capsule network model, Keras and tensor flow deep learning libraries with CUDA and CUDNN libraries for GPU expedited training. The model is trained using Google Colab 's framework. We may conclude that the capsule network is superior to the other methods described and performs better.

V . CONCLUSION

Detection of traffic signs is a difficult task and capsule networks use their inherent ability to detect the location and perform better compared to CNN's and capsule networks to improve reliability and accuracy by performing image detection and recognition tasks correctly, including on blurred, rotated and distorted images.

VI . REFERENCES

- [1] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recognition benchmark: a multi-class classification competition. In Neural Networks (IJCNN), The 2011 International Joint Conference on, pages 1453–1460. IEEE, 2011.
- [2] Geoffrey Hinton, Nicholas Frosst, and Sara Sabour. Matrix capsules with em routing. 2018.
- [3] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. CoRR, abs/1710.08864, 2017.
- [4] N. Deepika and V. V. Sajith Variyar. Obstacle classification and detection for vision based navigation for autonomous driving. In 2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI), pages 2092–2097, Sept 2017.
- [5] Yihui Wu, Yulong Liu, Jianmin Li, Huaping Liu, and Xiaolin Hu. Traffic sign detection based on convolutional neural networks. In Neural Networks (IJCNN), The 2013 International Joint Conference on, pages 1–7. IEEE, 2013.
- [6] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial Transformer Networks. ArXiv e-prints, June 2015.
- [7] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), volume 1, pages 886–893 vol. 1, June 2005.
- [8] David G. Lowe. Object recognition from local scale-invariant features. In Proceedings of the International Conference on

Computer Vision Volume 2 - Volume 2, ICCV '99, pages 1150–, Washington, DC, USA, 1999. IEEE Computer Society.

[9] Jung-Guk Park and Kyung-Joong Kim. Design of a visual perception model with edge-adaptive gabor filter and support vector machine for traffic sign detection. *Expert Systems with Applications*, 40(9):3679 – 3687, 2013.

[10] Rabia Malik, Javaid Khurshid, and Sana Nazir Ahmad. Road sign detection and recognition using colour segmentation, shape analysis and template matching. In *Machine Learning and Cybernetics, 2007 International Conference on*, volume 6, pages 3556–3560. IEEE, 2007.

[11] Fatin Zaklouta and Bogdan Stanciulescu. Real-time traffic sign recognition in three stages. *Robotics and autonomous systems*, 62(1):16–24, 2014.

[12] Andreas Mogelmose, Mohan Manubhai Trivedi, and Thomas B Moeslund. Vision-based traffic sign detection and analysis for intelligent driver assistance systems: Perspectives and survey. *IEEE Transactions on Intelligent Transportation Systems*, 13(4):1484–1497, 2012.

[13] C. K. Chandni, V. V. S. Variyar, and K. Guruvayurappan. Vision based closed loop pid controller design and implementation for autonomous car. In *2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pages 1928–1933, Sept 2017.

[14] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks*, 32:323–332, 2012.

[15] Pierre Sermanet and Yann LeCun. Traffic sign recognition with multiscale convolutional networks. In *IJCNN*, pages 2809–2813. IEEE.