

Vehicle Accident Detection

Alok Kumar, Indreesh Pandey, Arun Prakash, Akash Srivastava

Sinhgad College of Engineering, Pune, Maharashtra, India

Abstract: On Road accidents have become so rife that the fatalities due to these accidents have seen a steep rise. Accidents are a menace not only to human life, but resources, and it creates disruption in the normal traffic flow. The more concerning characteristic of the on road accidents is the overall damage rate that ensues these accidents. Now, when Machine Learning has taken over, the previously complex problems have become feasible, and the real life applications of these artificial ML models have been very promising. In this paper, we propose a learning Model that learns over an image dataset, thereby classifying never before seen images and data. With this model, we aim to predict and detect real-time collisions and accidents through a security camera planted along the periphery of the highways. We aim to classify these real-time accidents based on the level of damage. We will make use of the Artificial Neural Network to train the model to learn the similarities among images and accident data.

Keywords: Road accidents

I. INTRODUCTION

Artificial Neural Network is a computational system, wherein a number of layers of neurons are associated together, or are stacked one over the other in such a way that the output of one layer subsequently becomes the input of the next layer, until the final layer associates the input to a desired output. The prediction may not always be as expected, and in such a case, we apply a correction function that minimizes the error. This error correction is achieved using a back-propagation technique.

Neural networks use basis functions that follow the same form as, so that each basis function is itself a nonlinear function of a linear combination of the inputs, where the coefficients in the linear combination are adaptive parameters. This leads to the basic neural network model, which can be described as a series of functional transformations.

The Linear Classification Models don't work for complex feature set and dataset. This is due to the non-linearity and complex structure of data such as images, videos, audio, etc. The features that need to be taken into account is not specified or there are simply too many features present in the data. We then use complex neural Models that emulate Human Like intelligence through activation function implemented on neurons. The system models the learning process on dataset by minimizing the overall error due to the incorrect prediction. System like the Neural Network work by combining various elements and features of the dataset and processing these combinations in a number of layers. While most of the models implement a similar procedure, some exception may exist. Artificial Neural Network could be categorized under the learning procedure they tend to adopt.

II. LITERATURE SURVEY

A. A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects (Zewen Li, Wenjie Yang, Shouheng Peng, Fan Liu, Member, IEEE) [1]

A number of CNNs have been proposed. As we slowly begin to understand the shortcomings with a particular model, we try to minimize the effect of those shortcomings, which in turn, gives rise to newer extensions to the already established methods, creating methods that eke out efficiency and complexity. The advantages of the convolutional neural network have been discussed at length in the above-mentioned Literature. Some of these ideas are related to down-sampling (or dimensionality reduction), weight sharing and local connection. The paper gives the detailed survey of the CNN including applications, classic networks, building blocks, related functions and prospects. Construction of a CNN-based model from scratch has been taken up and discussed quite efficiently. Furthermore, some of the guidelines in devising a network aiming to improve speed and accuracy, as we know that building a deep neural network can improve the ability of the neural network to handle the complex tasks, such as low rank approximation, dimension reduction, faster computing speed. Certain optimizers have been discussed which give us an idea of selecting the appropriate functions. Various loss and activation functions have been discussed in a very alluring way, mentioning the advantages and shortcomings of every function that is used. The activation functions include sigmoid, tanh and ReLU. It also discusses certain other forms of the commonly-

used Non-linear Function ReLU such as PReLU and leaky ReLU. the paper also postulates various ways to choose and optimize error and loss functions. To name a few, MSE(Mean Square Error), MAE(Mean Absolute Error), Cross Entropy, Contrastive Error. Back- propagation has also been discussed sufficiently. In this paper, we saw some typical application of the CNN and the different dimensions used for the various problems. Also this paper gives us the challenges faced by the CNN such as - poor crowded scene result, lack of equivariance, low generalization, which makes CNN hard to handle.

B. Accident Forecasting in CCTV Traffic Camera Videos (Ankit Shah, Jean Baptiste Lamare, Tuan Nguyen Anh, Alexander Hauptmann) [2]

This paper has proposed the dataset CADP - Car Accident Detection and Prediction (CADP) dataset for Traffic Camera based accident forecasting. The discussion begins with the need for a new dataset, and thereby it establishes the basis of a new dataset. The advantages and limitations of other datasets are taken up at a greater length. The paper gives a detailed description about the CADP dataset right from the development to the results and its future scopes for application in different domains. Later on, the authors have also shown the results of classification using Faster R-CNN along with context mining and augmented context mining. Comparison is made between SSD and Faster R-CNN based upon experimental results on a new trainval set especially sampled for cross validation. They presented the results of state-of-the-art object detection and accident forecasting models on their dataset, highlighted the strengths and weaknesses of baseline models, and outperformed the initial results by adding context mining or augmented context mining. The results showed that augmented context mining does not improve the score obtained with a gradual context mining for object detection.

C. A Vision-Based Video Crash Detection Framework for Mixed Traffic Flow Environment Considering Low-Visibility Condition (Chen Wang ,Yulu Dai ,Wei Zhou ,and Yifei Geng) [3]

There has been various research conducted with the aim of developing vision-based system to detect traffic crashes. The paper mentions three categories of study:- (1) modelling the traffic flow (2) modelling of vehicle interactions, and (3) Analysis of activities. Analysing patterns of vehicle flow from previously seen data and models, it becomes important that real-time is consistent with the flow model, otherwise, the classification yields unsatisfactory results. Another method is to detect speed change through high sensitivity cameras. The third category mainly depends on using tracker facilities to monitor features such as speed, Force, acceleration, etc. Another practical limitation of the model is the environmental effects that hampers the efficiency of the model (foggy conditions, Torrential downpour and insufficient illumination). Various Image enhancing features could mitigate these effects. But, then, the complexity of the image-processing model increases manifold. The camera-based detection suffers mainly due to the computational complexity of the model, which brings down the latency in real-time implementation. Choosing a good learning and classification algorithm is also the key to a better Model.

D. Faster CNNs

- **R-CNN** : It uses a selective-search algorithm which extracts 2000 regions from an image and they are called region proposals. As we have to classify 2000 region proposals per image while training to network it will take a huge amount of time. Thus, it can not be implemented as it takes approx 47 seconds for each test image.
- **Fast R-CNN** : It uses a similar approach as the R- CNN. But, here we feed the total input image to CNN to generate convolutional feature maps instead of feeding the regional proposals. The Fast R-CNN is faster than R- CNN in case of training and testing as it does not feed 2000 regional proposals to CNN per image. But Including regional proposals decrease the efficiency of Fast R-CNN during testing time.
- **Faster R-CNN**: It allows the network learn region proposals. It feeds input images to CNN which provide convolutional Feature maps. Faster R-CNN uses separate networks to predict regional proposals instead of using selective-search algorithms. From the below graph, we can see that Faster R-CNN is much faster than its predecessors.

III. ARCHITECTURE REVIEW OF CNN

CNN is a deep neural network used in visual image processing. CNNs have managed to achieve superhuman performance on some complex visual tasks. They power image search services, self-driving cars, automatic video

classification systems, and more. Moreover, CNNs are not restricted to visual perception: they are also successful at other tasks, such as voice recognition or natural language processing (NLP); however, we will focus on visual applications of image recognition.

A. Convolutional Layer

The convolutional layer of a convolutional network is the most important part, which does most of the computational heavy lifting. The convolutional layer's parameters consist of a set of learnable filters. Every filter is small spatially (width and height), but extends through the full input volume. We slide the filter across width and height of the input volume and compute dot product between filter indices and corresponding input volume locations. As we slide the filter over the volume, we get a 2-d feature map. The network will learn filters that activate when they see some type of visual feature such as an edge of some orientation or a blotch of some color on the first layer. We stack the different feature maps obtained to get the output volume.

B. Pooling Layer

It is common to periodically insert a Pooling layer in-between successive Conv layers in a ConvNet architecture. Its function is to progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting. The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, using the MAX operation. The most common form is a pooling layer with filters of size 2x2 applied with a stride of 2 downsamples every depth slice in the input by 2 along both width and height, discarding 75% of the activations. Every MAX operation would in this case be taking a max over 4 numbers (little 2x2 region in some depth slice).

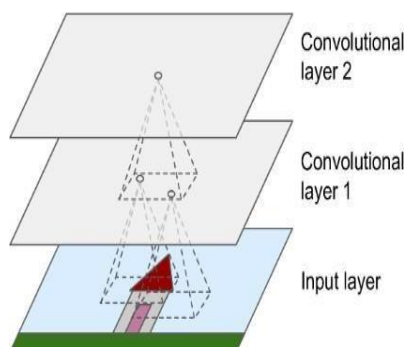


Fig. 1. Convolutional Layer

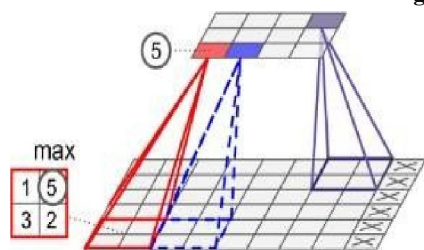


Fig. 2. Pooling Layer

CNN

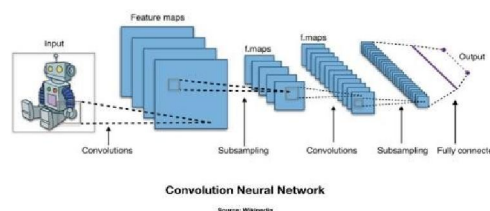


Fig. 3. Network Layers.

C. Fully-connected layer

Neurons in a fully connected layer have full connections to all activations in the previous layer, as seen in regular Neural Networks. Their activations can hence be computed with a matrix multiplication followed by a bias offset.

IV. DATASET

A. Accident Detection and Prediction (CADP) Dataset

The dataset consists of 1,416 video segments collected from YouTube, with 205 video segments having full spatio-temporal annotations. This dataset is largest in terms of number of traffic accidents, compared to other contemporary datasets. CADP contains videos collected from YouTube which are captured under various camera types and qualities, weather conditions and edited/resampled videos.

B. DETRAC Dataset

The UA-DETRAC benchmark dataset consists of 100 challenging videos captured from real-world traffic scenes (over 140,000 frames with rich annotations, including illumination, vehicle type, occlusion, truncation ratio, and vehicle bounding boxes) for multi-object detection and tracking.

C. Other Related Datasets

- The MIT dataset for traffic camera events
- NGSIM dataset for road traffic modeling
- CBSR dataset for single views at complex intersections
- CVRR dataset which simulated videos generated for traffic modeling
- Dashcam Accident Dataset (DAD)

V. TRAINING METHODOLOGY

We divide the training phase into a number of sub-phases, wherein the learning process is carried sequentially and different sub-phases are all dependent on the previous instances. We take a brief look into the proposed methodology.

A. DATASET PREPARATION

We use two datasets, namely CADP and DETRAC. CADP has accident images from various CCTV camera footage on youtube all accumulated at an single location. DETRAC, on the other hand, comprises non-accident, normal images that are derived from traffic cameras. We use 28000 accident images from CADP and 23000 non-accident images from DETRAC dataset. The images are divided into directories 'Accident' and 'Non-Accident', and each image is subsequently scaled to 128×128 and then converted to grayscale. The scaled vector with the image matrix and label are appended into a numpy array. After loading the complete set of images, we shuffle the contents of the array(shuffling inserts images with two different labels at random locations, which bring a randomness into the data). The image matrices and labels are store into two separate files on the disk.

Initially, the Dataset was generated from videos that consisted only a small fraction of positive images and large number of false positives. We call this Dataset a 'Crude Dataset'. This 'Crude Dataset' may produce poor validation results and generalization. The efficiency on the 'Crude Dataset' was mediocre.

In order to overcome the false positives, We manually validated the positive and false positive images, thereby removing most of the false positives. The Dataset was now reduced to 4000 Accident images. We took twice the number of 'Non-accident' Images, obtaining a sample size of 12000 images with relative ratio of 1:2 between 'Accident' and 'Non-Accident' images.

B. LOADING TRAINING DATA

The previously loaded image matrices and labels are retrieved for the learning procedure. We load the previously stored data from the disk. For each image in the loaded set, we nor- malize the intensity values at spatial locations. Normalization removes bias from dataset by downsizing the large contribution from some features and upsizing small contribution from some features.

C. NETWORK TRAINING

We have sequentially stacks all the layers - keras Conv2D, MaxPooling2D, Activation, Flatten, Dense, Dropout layer. The convolutional layer extracts the feature map from image vol- ume. Activation layer adds the activation function for

adding non-linearities. Flatten layer converts all output layers to 1D vectors. 'ReLu' and 'Sigmoid' are used as activation functions, and 'binary crossentropy' as our loss function and 'Adam' optimizer and for accuracy matrix calculation. We take three instances - 1000, 5000 and 10000 images for training the model on 'Crude Dataset'. We then evaluate the model on the newly created dataset of 12000 images. The whole sample was used for training.

For Testing, from the three above-mentioned instances, we take 10% instances as the validation set. We observe that for low-size training samples on 'Crude Dataset', the accuracy is low. On increasing, the number of instances, the accuracy somewhat increases, but, on further addition to the sample size, the accuracy hardly increases.

For the new Dataset of 12000 images, Testing was done on 350 images sample consisting of 'Accident' and 'Non-accident' images in equal proportions.

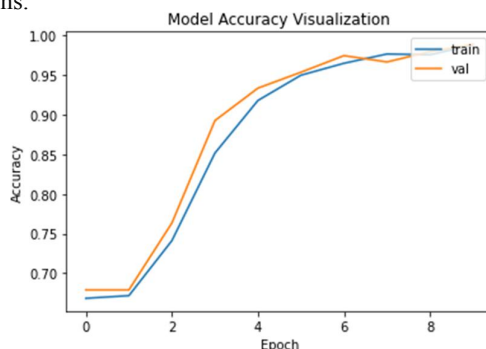


Fig. 6. Validation and Testing Accuracy

VI. PRELIMINARY RESULTS

1. The efficacy obtained so far is moderate. This maybe due to the problem of object detection, the small similarity between accident images, and the lack of more non-linear layers that learns the non-linear features.
2. YOLO(You look only once) Algorithm can be added to classify the various type of objects in the accident and non-accident images. This will enhance the accuracy of the model.
3. Context Augmentation with faster RCNN could further improve the accuracy.
4. Better accuracy can be achieved by optimizing the CADP dataset which contains false accident images

VII. EXPECTED RESULTS

After manual validation on 'Crude Dataset', we obtain a new 'Compact Dataset'. The result on this 'Compact Dataset' are provided below:

1. The Training and Validation Efficiency obtained was 97% and Testing Accuracy was 68% respectively.
2. The Recall and Precision were 1 and 0.68 respectively.
3. Cosine Similarity was recorded as 95%.

The various error and loss functions are included in the table given below:

Error and Loss Functions			
	Training	Validation	Testing
MSE	0.25	0.25	0.25
MAE	0.5	0.5	0.5
MAPE	164072304.00	160468032.00	2042376.00
Precision	0.67	0.68	0.59
Recall	0.99	1.00	1.00
Cosine S.	0.95	0.96	0.83

Fig. 7. Loss Function and Error Measure.

Volume 2, Issue 5, May 2022

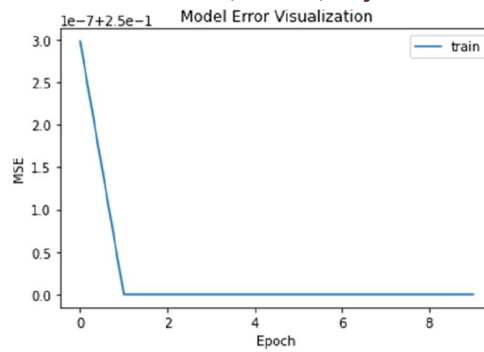


Fig. 8. Mean Squared Error

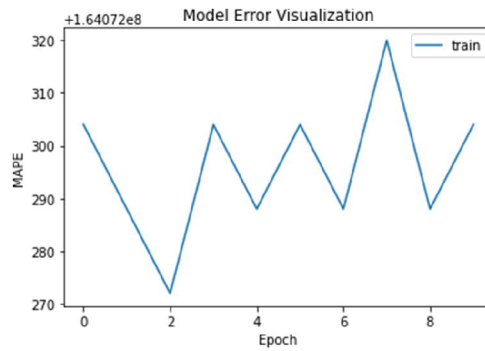


Fig. 9. Mean Absolute Percentage Error

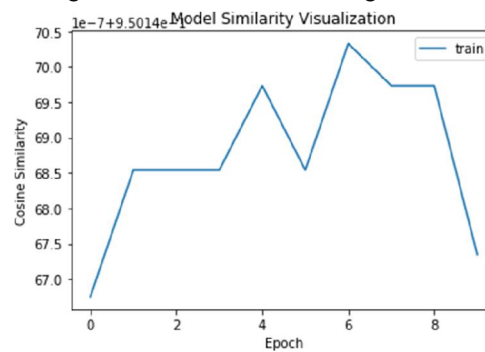


Fig. 10. Cosine Similarity

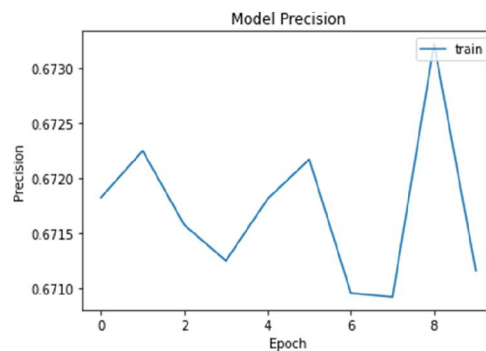


Fig. 11. Precision and Recall.

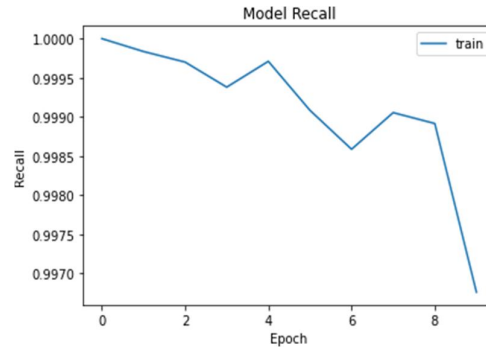


Fig. 12. Precision

REFERENCES

- [1]. Li, Zewen, Wenjie Yang, Shouheng Peng, and Fan Liu. "A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects." arXiv preprint arXiv:2004.02806 (2020).
- [2]. Shah, Ankit, Jean Baptiste Lamare, Tuan Nguyen Anh, and Alexander Hauptmann. "Accident forecasting in CCTV traffic camera videos." arXiv preprint arXiv:1809.05782 (2018).
- [3]. Wang, Chen, Yulu Dai, Wei Zhou, and Yifei Geng. "A vision-based video crash detection framework for mixed traffic flow environment considering low-visibility condition." Journal of advanced transportation 2020 (2020).
- [4]. Shah, Ankit Parag, Jean-Bapstite Lamare, Tuan Nguyen-Anh, and Alexander Hauptmann. "Cadp: A novel dataset for cctv traffic camera based accident analysis." In 2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), pp. 1-9.