

OBJECT DETECTION USING CONVOLUTION NEURAL NETWORK

N. K. Darwante

Assistant Professor, Department of Electronics & Tele.
Sanjivani College of Engineering Kopergaon, India, Affiliated to Savitribai Phule.
darwante11@gmail.com

Dr. U. B. Shinde

Dean Faculty of Engineering & Technology, Dr. B.A.M. University, Aurangabad, India.
shindeulhas1@yahoo.co.in

ABSTRACT

In Computer vision the object detection is one of the most important challenging points. The mostly best-performing detectors are based on the technique of convolution neural network. We investigate more than one techniques to improve the convolutional neural network based image recognition, distribution pipeline or detection of objects. In convolution neural network based technique of finding or searching region or location proposals in order to localize instances or objects. These activated regions or locations are mostly used as input to convolution neural network to express soulful features. In object detection the task of assigning a label and a bounding box to all instances or objects. We present a simple, easy yet strong formulation of object recognition or detection as a regression problem to object bounding box. And define a multi-scale inference procedure which is able to produce object detections with high resolution at a low cost by a network applications. The methods based on deep learning (Convolution neural network) have become the state of the art in object detection in image. State-of-the-art performance of the approach is shown on PASCAL VOC database. At last, we evaluate the mean average precision (mAP) for each class. The mAP of our model on PASCAL test dataset is 37.38%.

I.INTRODUCTION

The convolutional neural network(CNN, or ConvNet) is a class of deep, soulful, feed-forward artificial neural network that have successfully been applied to analyzing visual imagery. The receptive fields of various neurons partly overlapping so that they cover the entire visual field. The mostly important points in our projects are object detection which has tremendous, exorbitant application in our daily life. The main aim of this object detection is endorse more than one object in an only one image that is single images and to return the assurance of the class for each and every object and also predict the assuring bounding boxes. Due to its highest capacity in accurately classifying images CNNs tremendously useful in visual recognition since 2012 [1]. In [1], the authors show an extremely increment on the precision of image distribution in ILSVRC ImageNet Large Scale Visual Recognition Challenge.

In addition image distribution or classification, researchers have elongated the application of Convolution Neural Network to various other tasks in visual recognition that are segmentation [3], and object detection [5], generating sentences from image [4], localization [2],.

To sneeze about thousands of objects from millions of images, we have necessary a model with a largest learning capacity. However, the vast complexity of the object recognition task means that this problem cannot be specific even by a dataset as far as ImageNet. So our model should also have knowledge to recompense for all the data we don't have.

Besides, the presented method is quite simple. There is not necessary to hand design a model which captures or possession parts and their relations manifestly. This ingenuousness has the advantage of simple applicability to splay range of classes, but also show highest good detection performance beyond a higher range of objects – rigid ones and also deformable ones. This is presented together with state-of-the-art means the highest degree of development of an art or technique at particular time detection results on Pascal VOC challenge database [7].

Convolutional neural networks establish the class of models [11, 12, 13, 14]. This network capability can be possession by different their breadth, depth, and also make large strong and mainly accurate presumption about the image nature namely, stationary of statistics and locality of pixel dependencies. Therefore , comparison to standard networks with similarly-sized layers, CNNs have much small number of connections and also parameters and so they are simplifier to train, while their theoretically-better performance.

II.MOTIVATION AND OBJECTIVES

The CNN is a class of deep, artificial, feed forward neural network that have successfully been applied to analysing visual imagery. The attractive quality of CNN over other classification methods is that other methods are relatively inefficient with their local architecture and they are prohibitively expensive to apply for large scale high resolution images. The immense complexities of the object recognition task are easily reduced using CNN. The CNN useful when multiple objects in image and objects are small and detect exact location and size of object in image is desired.

By using convolution neural network to classify the millions of high resolution Images are classified in ImageNet ISVRC like contest into the thousands of different classes. The NMS (non maximal suppression) technique is used to merge highly overlapping region which is predicted to be of same class.

In SVM (Support vector machine) and bounding box regression combined together to provide a high performance in object detection. The CNN models have large learning capacity.The CNN model the mostly helpful in video recognition and more application that are localization, segmentation, generating sentence from image and object detection and recognitions.

III.BACKGROUND

Object detection is the most attractive point in visual recognition area in the past years.Today people tend to design features from raw images to improvement in the performance of the detection. These features SIFT [9] and HOG [10] are most successful and these features are combined with SVMs have detect pedestrian from images. Detection and classification of object are fundamental building blocks of artificial intelligenceand detection of objects refers to the determination of the presence or absence of specific features in image data. These specific features are detected; an object can be further classified as belonging to one of a pre-defined set of classes. This latter operation is also called as object classification.

The people focus on CNN since[1] it shows improvement and accuracy and efficiency in classification. By using deep CNN and to localization of the object by specifying a bounding edge box therefore we cannot directly used CNN in direction of object because solving the localization problem [2].All round overall nature of sliding window with high-computational complexity. This method cannot be practically implemented and its approach cannot works in practice.

Sliding Windows method is to provide an idea for solving localization by classifying regional proposals of the image. In the R-CNN [5] the researcher choose selective search as proposal generation algorithm due to its fast computational times.

Many people also develop supervised learning and deformable CNNs [12] to detect objects in past years. Now a day's very vast techniques are used in detection and recognition of objects.

IV.RELETED WORK

In this system we develop a new approach for detecting multiple objects from images based on convolutional neural networks (CNNs). Mostly works in the detection of objects, regional CNNs(rCNN) [5] is the important remarkable one that combines with selective search[6], CNNs, support vector machines(SVM) and bounding box regression together and to provide a higher performance in detection and recognition of object.

R-CNNs for object detection were first presented in 2014 by Ross girshick et al, and state-of-the-art approaches on one of the major objects recognition challenges in the field: Pascal VOC. The basic idea of R-CNN is to take a deep neural network which was originally trained for image classification using millions of annotated images and modify it for the purpose of object detection.

Selective search: Selective search is a method for finding a large set of possible object Locations in an image, independent of the class of the actual object. It works by clustering image pixels into segments, and then performing hierarchical clustering to combine segments from the same object into object proposal.

NMS (Non Maximum Suppression):

Object detection methods often output multiple detections which fully or partly cover the same object in an image. These ROIs (Regions-of-interests) need to be merged to be able to count objects and obtain their exact locations in the image. This is traditionally done using a technique called Non Maximum Suppression (NMS).

The version of NMS we use (and which was also used in the R-CNN publications) does not merge ROIs but instead tries to identify which ROIs best cover the real locations of an object and discards all other ROIs. This is implemented by iteratively selecting the ROI with highest confidence and removing all other ROIs which significantly overlap this ROI and are classified to be of the same class.

mAP (mean Average Precision):

Once trained, the quality of the model can be measured using different criteria, such as precision, recall, accuracy, area-under-curve, etc. A common metric which is used for the Pascal VOC object recognition challenge is to measure the Average Precision (AP) for each class. The following description of Average Precision is taken from Everinghamet.al. The mean Average Precision (mAP) is computed by taking the average over the APs of all classes.

For a given task and class, the precision/recall curve is computed from a method's ranked output. Recall is defined as the proportion of all positive examples ranked above a given rank. Precision is the proportion of all examples above that rank which are from the positive class. The AP summarizes the shape of the precision/recall curve, and is defined as the mean precision at a set of eleven equally spaced recall levels $[0, 0.1, \dots, 1]$:

$$AP = \frac{1}{11} \sum_{r \in \{0, 0.1, \dots, 1\}} p_{interp}(r)$$

The precision at each recall level r is interpolated by taking the maximum precision measured for a method for which the corresponding recall exceeds r :

$$p_{interp}(r) = \max_{\tilde{r}: \tilde{r} \geq r} p(\tilde{r})$$

Where $p(\tilde{r})$ is the measured precision at recall \tilde{r} . The intention in interpolating the precision/recall curve in this way is to reduce the impact of the “wiggles” in the precision/recall curve, caused by small variations in the ranking of examples. It should be noted that to obtain a high score, a method must have precision at all levels of recall – this penalizes methods which retrieve only a subset of examples with high precision (e.g. side views of cars).

V.OVERVIEW OF SYSTEM

In the system the class-agnostic scalable object detection achieving by predicting a set of bounding boxes this is important points which represents the potential objects. The edge boxes algorithm uses as our proposal generation algorithm. This algorithm generates and scores the proposal based on the edge map of the image is the basic idea of edges boxes. This system generates an edge map with a structured edge detector where each pixel contains a NMS (non-maximal suppression) to merge or join the bounding box to get the proposals and we evaluate the mean average precision (mAP) for each class. The mAP of our model on PASCAL test dataset is 37.38%. The selective search which has been chosen in the rCNN, [8] shows that for VOC dataset, the mAP of edge boxes is which is slightly larger than selective search with mAP as 31.7%. edge boxes is that the runtime is much faster than the majority of the proposal generation schemes. For edge boxes, the average runtime is 0.3 seconds while for selective search as 10 seconds. Hence the edge boxes decreases the time complexity without degrading the performance. Therefore, mostly choose the edge boxes as the proposal generation algorithm.

VI. CONCLUSION

In this work, we developed a deep learning method for solving the edge detection problem by using convolutional neural networks (CNN).

We provide a new model for detection of objects based on Convolution neural network. And this model we use the edge boxes algorithm to generate or produce new proposals, and uses of a fine-tuned the CaffeNet model to produce or generate the score for each proposals.

In previous work, our approach does not need extra additional feature extraction process and it can be very simple, easy and fast forward while achieving better result. And also very easy simple for people to implement and integrate our easy algorithm into their own computer vision systems. The deep learning becoming more and more popular, it is probable that in the network can be combined in other of the network used in application.

And we will changes in the deeper network to increase the accuracy, clarity of classification and also to add the ground truth bounding boxes into the training data to improve the localization clarity and accuracy of images.

VII. ACKNOWLEDGMENT

We take this opportunity to thank all the people who have given assistance to us in the form of advice, suggestions, and any other for completion of this paper. It is a pleasure to convey our gratitude to them all in our humble acknowledgment. It is an honour for us to express our sincere gratitude to Dr. D. N. Kyatanavar, Principal Sanjivani College of Engineering, and Kopargaon for providing all necessary support for completion of this paper. Specifically, thanks to Head of Department of electronics and telecommunication, Dr. B. S. Agarkar, for the continuous support, motivation, enthusiasm, and immense knowledge. Next, it is our pleasure to acknowledge the various individuals, who contributed in completion of our paper. Finally, we are thankful to our family for allowing us to complete this paper in the time.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In NIPS, 2012.
- [2] C. Szegedy, A. Toshev, and D. Erhan. Deep neural networks for object detection. In NIPS, 2013. 2.
- [3] L. Jonathan, S. Evan, D. Trevor, Fully convolutional networks for semantic segmentation, to appear in CVPR 2015.
- [4] K. Andrej, Li Fei-Fei, Deep visual-semantic alignments for generating image descriptions. arXiv:1412.2306.
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. in Proc. CVPR, 2014.
- [6] J. Uijlings, K. vande Sande, T. Gevers and A. Smeulders, Selective search for object recognition. in Proc. IJCV, 2013.
- [7] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal visual object classes (voc) challenge. International Journal of Computer Vision, 88(2):303–338, 2010.
- [8] H. Jan, B. Rodrigo, S. Bernt. How good are detection proposals, really? arXiv:1406.6962.
- [9] R. Girshick, J. Donahue, T. Darrell and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. in Proc. CVPR, 2014.
- [10] Y. Le Cun, B. Boser, J.S. Denker, D. Henderson, R.E. Howard, W. Hubbard, L.D. Jackel, et al. Handwritten digit recognition with a back propagation network. In Advances in neural information processing systems, 1990.
- [11] K. Jarrett, K. Kavukcuoglu, M. A. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In International Conference on Computer Vision, pages 2146–2153. IEEE, 2009.
- [18] H. Lee, R. Grosse, R. Ranganath, and A.Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In Proceedings of the 26th Annual International Conference on Machine Learning, pages 609–616. ACM, 2009.
- [16] Y. LeCun, F.J. Huang, and L. Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on, volume 2, pages II–97. IEEE, 2004.
- [13] N. Pinto, D. Doukhan, J.J. DiCarlo, and D.D. Cox. A high-throughput screening approach to discovering good forms of biologically inspired visual representation. PLoS computational biology, 5(11):e1000579, 2009.