

Network In Network

Min Lin^{1,2}, Qiang Chen², Shuicheng Yan²

¹Graduate School for Integrative Sciences and Engineering

²Department of Electronic & Computer Engineering

National University of Singapore, Singapore

{linmin, chenqiang, eleyans}@nus.edu.sg

Abstract

We propose a novel deep network structure called “Network In Network”(NIN) to enhance model discriminability for local patches within the receptive field. The conventional convolutional layer uses linear filters followed by a nonlinear activation function to scan the input. Instead, we build micro neural networks with more complex structures to abstract the data within the receptive field. We instantiate the micro neural network with a multilayer perceptron, which is a potent function approximator. The feature maps are obtained by sliding the micro networks over the input in a similar manner as CNN; they are then fed into the next layer. Deep NIN can be implemented by stacking multiple of the above described structure. With enhanced local modeling via the micro network, we are able to utilize global average pooling over feature maps in the classification layer, which is easier to interpret and less prone to overfitting than traditional fully connected layers. We demonstrated the state-of-the-art classification performances with NIN on CIFAR-10 and CIFAR-100, and reasonable performances on SVHN and MNIST datasets.

1 Introduction

Vấn đề của mô hình tuyến tính tổng quát (GLM):

- Bộ lọc tích chập trong CNN được xem như một mô hình tuyến tính tổng quát (GLM) trên các dữ liệu cục bộ, tuy nhiên mức độ trừu tượng của mô hình này bị cho là thấp. Sự trừu tượng ở đây ám chỉ khả năng mô hình có thể nhận biết các biến thể của cùng một khái niệm (concept).
- GLM chỉ hoạt động tốt khi các biến thể của khái niệm này có thể phân tách tuyến tính, tức là các biến thể nằm cùng một phía của mặt phẳng phân tách được xác định bởi GLM.

Giải pháp thay thế GLM:

- Tác giả đề xuất thay thế GLM bằng một bộ xấp xỉ hàm phi tuyến mạnh hơn. Trong trường hợp này, tác giả sử dụng cấu trúc “mạng vi mô” (micro network) thay thế, với mạng MLP (multilayer perceptron). MLP là bộ xấp xỉ hàm phi tuyến tổng quát và có khả năng huấn luyện bằng phương pháp lan truyền ngược (back-propagation).

Cấu trúc của Network in Network (NIN):

- NIN sử dụng nhiều lớp mlpconv, là sự kết hợp của mạng vi mô với MLP. Thay vì chỉ đơn giản sử dụng các lớp tích chập như CNN truyền thống, mlpconv tạo ra một mạng cục bộ vi mô để ánh xạ vùng thụ cảm cục bộ của đầu vào thành vector đặc trưng đầu ra. Điều này tạo ra khả năng mô hình hóa cục bộ mạnh mẽ hơn.

Pooling toàn cục (Global Average Pooling):

- Thay vì sử dụng các lớp fully connected truyền thống, tác giả sử dụng phép pooling toàn cục để tính trung bình các bản đồ đặc trưng từ lớp mlpconv cuối cùng và đưa ra xác suất của từng loại.
- Ưu điểm của việc này là phép pooling toàn cục giúp mô hình dễ giải thích hơn và giảm thiểu vấn đề quá khớp (overfitting) mà các lớp fully connected có thể gặp phải. Đồng thời, nó cũng có vai trò như một bộ điều chỉnh (regularizer) giúp ngăn chặn quá khớp.

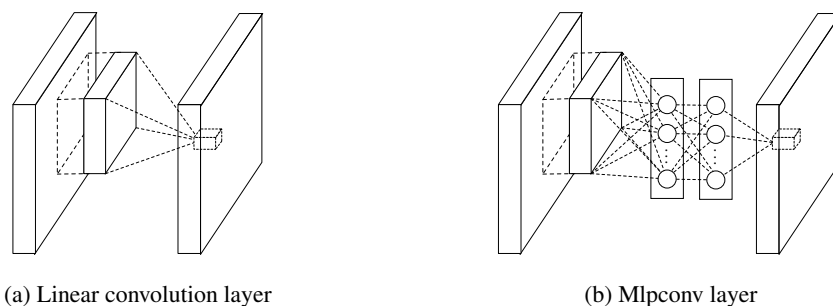


Figure 1: Comparison of linear convolution layer and mlpconv layer. The linear convolution layer includes a linear filter while the mlpconv layer includes a micro network (we choose the multilayer perceptron in this paper). Both layers map the local receptive field to a confidence value of the latent concept.

over the input in a similar manner as CNN and are then fed into the next layer. The overall structure of the NIN is the stacking of multiple mlpconv layers. It is called “Network In Network” (NIN) as we have micro networks (MLP), which are composing elements of the overall deep network, within mlpconv layers,

Instead of adopting the traditional fully connected layers for classification in CNN, we directly output the spatial average of the feature maps from the last mlpconv layer as the confidence of categories via a global average pooling layer, and then the resulting vector is fed into the softmax layer. In traditional CNN, it is difficult to interpret how the category level information from the objective cost layer is passed back to the previous convolution layer due to the fully connected layers which act as a black box in between. In contrast, global average pooling is more meaningful and interpretable as it enforces correspondance between feature maps and categories, which is made possible by a stronger local modeling using the micro network. Furthermore, the fully connected layers are prone to overfitting and heavily depend on dropout regularization [4] [5], while global average pooling is itself a structural regularizer, which natively prevents overfitting for the overall structure.

2 Convolutional Neural Networks

Classic convolutional neuron networks [1] consist of alternatively stacked convolutional layers and spatial pooling layers. The convolutional layers generate feature maps by linear convolutional filters followed by nonlinear activation functions (rectifier, sigmoid, tanh, etc.). Using the linear rectifier as an example, the feature map can be calculated as follows:

$$f_{i,j,k} = \max(w_k^T x_{i,j}, 0). \quad (1)$$

Here (i, j) is the pixel index in the feature map, $x_{i,j}$ stands for the input patch centered at location (i, j) , and k is used to index the channels of the feature map.

This linear convolution is sufficient for abstraction when the instances of the latent concepts are linearly separable. However, representations that achieve good abstraction are generally highly non-linear functions of the input data. In conventional CNN, this might be compensated by utilizing an over-complete set of filters [6] to cover all variations of the latent concepts. Namely, individual linear filters can be learned to detect different variations of a same concept. However, having too many filters for a single concept imposes extra burden on the next layer, which needs to consider all combinations of variations from the previous layer [7]. As in CNN, filters from higher layers map to larger regions in the original input. It generates a higher level concept by combining the lower level concepts from the layer below. Therefore, we argue that it would be beneficial to do a better abstraction on each local patch, before combining them into higher level concepts.

In the recent maxout network [8], the number of feature maps is reduced by maximum pooling over affine feature maps (affine feature maps are the direct results from linear convolution without

1. CNN cơ bản và công thức tính toán:

- CNN được cấu tạo từ các lớp tích chập (convolutional layers) và lớp pooling. Lớp tích chập tạo ra các bản đồ đặc trưng (feature maps) bằng cách sử dụng các bộ lọc tích chập tuyến tính. Công thức tổng quát cho một phép tích chập là:

$$f_{i,j,k} = \max(w_k^T x_{i,j}, 0)$$

Trong đó, $f_{i,j,k}$ là giá trị pixel trên bản đồ đặc trưng tại vị trí (i, j) , $x_{i,j}$ là vùng tử âm của đầu vào, và k là chỉ số của kênh trên bản đồ đặc trưng.

- Hạn chế của CNN tuyến tính:**
 - CHEN với các bộ lọc tuyến tính chỉ có khả năng trừu tượng hóa các khái niệm khi các khái niệm tiềm ẩn có thể phân tách tuyến tính. Tuy nhiên, các biểu diễn phức tạp hơn thường không thể được giải quyết tốt bởi các phép tuyến tính.
 - Để giải quyết vấn đề này, một giải pháp thông thường là sử dụng nhiều bộ lọc hơn để bao quát tất cả các biến thể của khái niệm tiềm ẩn. Điều này có nghĩa là cần có nhiều bộ lọc cho cùng một khái niệm.
- Tác động của việc sử dụng quá nhiều bộ lọc:**
 - Tác giả nhận mạnh rằng việc sử dụng quá nhiều bộ lọc cho một khái niệm có thể tạo thêm gánh nặng cho các lớp tiếp theo. Cụ thể, các lớp sau đó sẽ cần phải lý tính cả các kết hợp của những biến thể từ lớp trước, dẫn đến sự phức tạp và khó khăn trong việc tổng hợp các khái niệm cấp cao hơn.
- Giải pháp cho tác giả:**
 - Tác giả đề xuất rằng cần phải có một sự trừu tượng hóa tốt hơn ngay tại mỗi vùng cục bộ trước khi kết hợp chúng thành các khái niệm cấp cao. Điều này sẽ giảm bớt gánh nặng cho các lớp tiếp theo và giúp mô hình có khả năng trừu tượng hóa tốt hơn.

- Giải pháp của tác giả:**
 - Tác giả đề xuất rằng cần phải có một sự trừu tượng hóa tốt hơn ngay tại mỗi vùng cục bộ trước khi kết hợp chúng thành các khái niệm cấp cao. Điều này sẽ giảm bớt gánh nặng cho các lớp tiếp theo và giúp mô hình có khả năng trừu tượng hóa tốt hơn.
- Sơ lược về mô hình Maxout Network:**
 - Tác giả cũng đề cập đến mô hình Maxout, một mạng tích chập khác có khả năng xấp xỉ các hàm phi tuyến bằng một hàm xấp xỉ từng đoạn tuyến tính. Mô hình này giúp giải quyết một số hạn chế của CNN tuyến tính bằng cách tách biệt các khái niệm trong các lớp hợp lệ (linear sets).
 - Tuy nhiên, Maxout Network đặt ra giả thiết rằng các khái niệm tiềm ẩn nằm trong một tập hợp lồi trong không gian đầu vào, điều này không phải lúc nào cũng đúng trong thực tế.
- Giải pháp đề xuất: Network In Network (NIN):**
 - Để giải quyết vấn đề này, tác giả giới thiệu NIN, trong đó một mạng vi mô (micro network) được tích hợp vào mỗi lớp tích chập để tạo ra các đặc trưng trừu tượng hơn cho từng vùng cục bộ.
 - Cấu trúc NIN bao gồm việc sử dụng các mạng con MLP (Multilayer Perceptron) để tạo ra khả năng trừu tượng hóa mạnh mẽ hơn, và sử dụng các phép hợp nhất mô hình trong việc phân tích các khái niệm phức tạp.

applying the activation function). Maximization over linear functions makes a piecewise linear approximator which is capable of approximating any convex functions. Compared to conventional convolutional layers which perform linear separation, the maxout network is more potent as it can separate concepts that lie within convex sets. This improvement endows the maxout network with the best performances on several benchmark datasets.

However, maxout network imposes the prior that instances of a latent concept lie within a convex set in the input space, which does not necessarily hold. It would be necessary to employ a more general function approximator when the distributions of the latent concepts are more complex. We seek to achieve this by introducing the novel “Network In Network” structure, in which a micro network is introduced within each convolutional layer to compute more abstract features for local patches.

Sliding a micro network over the input has been proposed in several previous works. For example, the Structured Multilayer Perceptron (SMLP) [9] applies a shared multilayer perceptron on different patches of the input image; in another work, a neural network based filter is trained for face detection [10]. However, they are both designed for specific problems and both contain only one layer of the sliding network structure. NIN is proposed from a more general perspective, the micro network is integrated into CNN structure in pursuit of better abstractions for all levels of features.

3 Network In Network

We first highlight the key components of our proposed “Network In Network” structure: the MLP convolutional layer and the global averaging pooling layer in Sec. 3.1 and Sec. 3.2 respectively. Then we detail the overall NIN in Sec. 3.3.

3.1 MLP Convolution Layers

Given no priors about the distributions of the latent concepts, it is desirable to use a universal function approximator for feature extraction of the local patches, as it is capable of approximating more abstract representations of the latent concepts. Radial basis network and multilayer perceptron are two well known universal function approximators. We choose multilayer perceptron in this work for two reasons. First, multilayer perceptron is compatible with the structure of convolutional neural networks, which is trained using back-propagation. Second, multilayer perceptron can be a deep model itself, which is consistent with the spirit of feature re-use [2]. This new type of layer is called mlpconv in this paper, in which MLP replaces the GLM to convolve over the input. Figure 1 illustrates the difference between linear convolutional layer and mlpconv layer. The calculation performed by mlpconv layer is shown as follows:

$$\begin{aligned} f_{i,j,k_1}^1 &= \max(w_{k_1}^1 T x_{i,j} + b_{k_1}, 0). \\ &\vdots \\ f_{i,j,k_n}^n &= \max(w_{k_n}^n T f_{i,j}^{n-1} + b_{k_n}, 0). \end{aligned} \quad (2)$$

Here n is the number of layers in the multilayer perceptron. Rectified linear unit is used as the activation function in the multilayer perceptron.

From cross channel (cross feature map) pooling point of view, Equation 2 is equivalent to cascaded cross channel parametric pooling on a normal convolution layer. Each pooling layer performs weighted linear recombination on the input feature maps, which then go through a rectifier linear unit. The cross channel pooled feature maps are cross channel pooled again and again in the next layers. This cascaded cross channel parametric pooling structure allows complex and learnable interactions of cross channel information.

The cross channel parametric pooling layer is also equivalent to a convolution layer with 1x1 convolution kernel. This interpretation makes it straightforward to understand the structure of NIN.



Figure 2: The overall structure of Network In Network. In this paper the NINs include the stacking of three mlpconv layers and one global average pooling layer.

Comparison to maxout layers: the maxout layers in the maxout network performs max pooling across multiple affine feature maps [8]. The feature maps of maxout layers are calculated as follows:

$$f_{i,j,k} = \max_m (w_{k_m}^T x_{i,j}). \quad (3)$$

Maxout over linear functions forms a piecewise linear function which is capable of modeling any convex function. For a convex function, samples with function values below a specific threshold form a convex set. Therefore, by approximating convex functions of the local patch, maxout has the capability of forming separation hyperplanes for concepts whose samples are within a convex set (i.e. l_2 balls, convex cones). Mlpconv layer differs from maxout layer in that the convex function approximator is replaced by a universal function approximator, which has greater capability in modeling various distributions of latent concepts.

3.2 Global Average Pooling

Conventional convolutional neural networks perform convolution in the lower layers of the network. For classification, the feature maps of the last convolutional layer are vectorized and fed into fully connected layers followed by a softmax logistic regression layer [4] [8] [11]. This structure bridges the convolutional structure with traditional neural network classifiers. It treats the convolutional layers as feature extractors, and the resulting feature is classified in a traditional way.

However, the fully connected layers are prone to overfitting, thus hampering the generalization ability of the overall network. Dropout is proposed by Hinton et al. [5] as a regularizer which randomly sets half of the activations to the fully connected layers to zero during training. It has improved the generalization ability and largely prevents overfitting [4].

In this paper, we propose another strategy called global average pooling to replace the traditional fully connected layers in CNN. The idea is to generate one feature map for each corresponding category of the classification task in the last mlpconv layer. Instead of adding fully connected layers on top of the feature maps, we take the average of each feature map, and the resulting vector is fed directly into the softmax layer. One advantage of global average pooling over the fully connected layers is that it is more native to the convolution structure by enforcing correspondences between feature maps and categories. Thus the feature maps can be easily interpreted as categories confidence maps. Another advantage is that there is no parameter to optimize in the global average pooling thus overfitting is avoided at this layer. Furthermore, global average pooling sums out the spatial information, thus it is more robust to spatial translations of the input.

We can see global average pooling as a structural regularizer that explicitly enforces feature maps to be confidence maps of concepts (categories). This is made possible by the mlpconv layers, as they makes better approximation to the confidence maps than GLMs.

3.3 Network In Network Structure

The overall structure of NIN is a stack of mlpconv layers, on top of which lie the global average pooling and the objective cost layer. Sub-sampling layers can be added in between the mlpconv

layers as in CNN and maxout networks. Figure 2 shows an NIN with three mlpconv layers. Within each mlpconv layer, there is a three-layer perceptron. The number of layers in both NIN and the micro networks is flexible and can be tuned for specific tasks.

4 Experiments

4.1 Overview

We evaluate NIN on four benchmark datasets: CIFAR-10 [12], CIFAR-100 [12], SVHN [13] and MNIST [1]. The networks used for the datasets all consist of three stacked mlpconv layers, and the mlpconv layers in all the experiments are followed by a spatial max pooling layer which down-samples the input image by a factor of two. As a regularizer, dropout is applied on the outputs of all but the last mlpconv layers. Unless stated specifically, all the networks used in the experiment section use global average pooling instead of fully connected layers at the top of the network. Another regularizer applied is weight decay as used by Krizhevsky et al. [4]. Figure 2 illustrates the overall

2. Cấu trúc mạng sử dụng:

- Kiến trúc mạng được sử dụng cho tất cả các bộ dữ liệu đều bao gồm ba lớp mlpconv xếp chồng lên nhau. Đây là điểm khác biệt chính của NIN so với CNN truyền thống.
- Sau mỗi lớp mlpconv, một lớp max pooling không gian được thêm vào để giảm kích thước của hình ảnh đầu vào theo tỷ lệ 1:2.
- Các kỹ thuật điều chuẩn (regularization):
 - Dropout được áp dụng cho các đầu ra của tất cả các lớp mlpconv ngoại trừ lớp mlpconv cuối cùng. Dropout là một kỹ thuật phổ biến trong deep learning để ngăn chặn hiện tượng quá khớp (overfitting).
 - Tác giả cũng sử dụng global average pooling thay vì các lớp fully connected truyền thống ở đầu ra của mạng. Điều này giúp mô hình trở nên dễ giải thích hơn và giảm thiểu quá khớp.
 - Weight decay (suy giảm trọng số) được áp dụng như một kỹ thuật điều chuẩn bổ sung, giúp giảm thiểu overfitting bằng cách phạt các trọng số có giá trị lớn trong quá trình huấn luyện.

h Cấu hình huấn luyện:

- Tác giả sử dụng mã nguồn cuda-convnet phát triển bởi Alex Krizhevsky, một tác giả nổi tiếng trong cộng đồng học sâu nhờ đóng góp vào sự phát triển của kiến trúc AlexNet.
- Quá trình huấn luyện được thực hiện với các mini-batch có kích thước 128.
- Trong quá trình huấn luyện, sau khi độ chính xác trên tập huấn luyện không còn cải thiện, tác giả giảm learning rate theo hệ số 10 và tiếp tục huấn luyện. Quá trình này được lặp lại một lần nữa cho đến khi learning rate cuối cùng giảm xuống còn 1% so với giá trị ban đầu.

4.2 CIFAR-10

The CIFAR-10 dataset [12] is composed of 10 classes of natural images with 50,000 training images in total, and 10,000 testing images. Each image is an RGB image of size 32x32. For this dataset, we apply the same global contrast normalization and ZCA whitening as was used by Goodfellow et al. in the maxout network [8]. We use the last 10,000 images of the training set as validation data.

The number of feature maps for each mlpconv layer in this experiment is set to the same number as in the corresponding maxout network. Two hyper-parameters are tuned using the validation set, i.e. the local receptive field size and the weight decay. After that the hyper-parameters are fixed and we re-train the network from scratch with both the training set and the validation set. The resulting model is used for testing. We obtain a test error of 10.41% on this dataset, which improves more than one percent compared to the state-of-the-art. A comparison with previous methods is shown in Table 1.

Table 1: Test set error rates for CIFAR-10 of various methods.

Method	Test Error
Stochastic Pooling [11]	15.13%
CNN + Spearmint [14]	14.98%
Conv. maxout + Dropout [8]	11.68%
NIN + Dropout	10.41%
CNN + Spearmint + Data Augmentation [14]	9.50%
Conv. maxout + Dropout + Data Augmentation [8]	9.38%
DropConnect + 12 networks + Data Augmentation [15]	9.32%
NIN + Dropout + Data Augmentation	8.81%

It turns out in our experiment that using dropout in between the mlpconv layers in NIN boosts the performance of the network by improving the generalization ability of the model. As is shown in Figure 3, introducing dropout layers in between the mlpconv layers reduced the test error by more than 20%. This observation is consistent with Goodfellow et al. [8]. Thus dropout is added

in between the mlpconv layers to all the models used in this paper. The model without dropout regularizer achieves an error rate of 14.51% for the CIFAR-10 dataset, which already surpasses many previous state-of-the-arts with regularizer (except maxout). Since performance of maxout without dropout is not available, only dropout regularized version are compared in this paper.



Figure 3: The regularization effect of dropout in between mlpconv layers. Training and testing error of NIN with and without dropout in the first 200 epochs of training is shown.

To be consistent with previous works, we also evaluate our method on the CIFAR-10 dataset with translation and horizontal flipping augmentation. We are able to achieve a test error of 8.81%, which sets the new state-of-the-art performance.

4.3 CIFAR-100

The CIFAR-100 dataset [12] is the same in size and format as the CIFAR-10 dataset, but it contains 100 classes. Thus the number of images in each class is only one tenth of the CIFAR-10 dataset. For CIFAR-100 we do not tune the hyper-parameters, but use the same setting as the CIFAR-10 dataset. The only difference is that the last mlpconv layer outputs 100 feature maps. A test error of 35.68% is obtained for CIFAR-100 which surpasses the current best performance without data augmentation by more than one percent. Details of the performance comparison are shown in Table 2.

Table 2: Test set error rates for CIFAR-100 of various methods.

Method	Test Error
Learned Pooling [16]	43.71%
Stochastic Pooling [11]	42.51%
Conv. maxout + Dropout [8]	38.57%
Tree based priors [17]	36.85%
NIN + Dropout	35.68%

4.4 Street View House Numbers

The SVHN dataset [13] is composed of 630,420 32x32 color images, divided into training set, testing set and an extra set. The task of this data set is to classify the digit located at the center of each image. The training and testing procedure follow Goodfellow et al. [8]. Namely 400 samples per class selected from the training set and 200 samples per class from the extra set are used for validation. The remainder of the training set and the extra set are used for training. The validation set is only used as a guidance for hyper-parameter selection, but never used for training the model.

Preprocessing of the dataset again follows Goodfellow et al. [8], which was a local contrast normalization. The structure and parameters used in SVHN are similar to those used for CIFAR-10, which consist of three mlpconv layers followed by global average pooling. For this dataset, we obtain a

Table 3: Test set error rates for SVHN of various methods.

Method	Test Error
Stochastic Pooling [11]	2.80%
Rectifier + Dropout [18]	2.78%
Rectifier + Dropout + Synthetic Translation [18]	2.68%
Conv. maxout + Dropout [8]	2.47%
NIN + Dropout	2.35%
Multi-digit Number Recognition [19]	2.16%
DropConnect [15]	1.94%

test error rate of 2.35%. We compare our result with methods that did not augment the data, and the comparison is shown in Table 3.

4.5 MNIST

The MNIST [1] dataset consists of hand written digits 0-9 which are 28x28 in size. There are 60,000 training images and 10,000 testing images in total. For this dataset, the same network structure as used for CIFAR-10 is adopted. But the numbers of feature maps generated from each mlpconv layer are reduced. Because MNIST is a simpler dataset compared with CIFAR-10; fewer parameters are needed. We test our method on this dataset without data augmentation. The result is compared with previous works that adopted convolutional structures, and are shown in Table 4.

Table 4: Test set error rates for MNIST of various methods.

Method	Test Error
2-Layer CNN + 2-Layer NN [11]	0.53%
Stochastic Pooling [11]	0.47%
NIN + Dropout	0.47%
Conv. maxout + Dropout [8]	0.45%

We achieve comparable but not better performance (0.47%) than the current best (0.45%) since MNIST has been tuned to a very low error rate.

4.6 Global Average Pooling as a Regularizer

Global average pooling layer is similar to the fully connected layer in that they both perform linear transformations of the vectorized feature maps. The difference lies in the transformation matrix. For global average pooling, the transformation matrix is prefixed and it is non-zero only on block diagonal elements which share the same value. Fully connected layers can have dense transformation matrices and the values are subject to back-propagation optimization. To study the regularization effect of global average pooling, we replace the global average pooling layer with a fully connected layer, while the other parts of the model remain the same. We evaluated this model with and without dropout before the fully connected linear layer. Both models are tested on the CIFAR-10 dataset, and a comparison of the performances is shown in Table 5.

Table 5: Global average pooling compared to fully connected layer.

Method	Testing Error
mlpconv + Fully Connected	11.59%
mlpconv + Fully Connected + Dropout	10.88%
mlpconv + Global Average Pooling	10.41%

As is shown in Table 5, the fully connected layer without dropout regularization gave the worst performance (11.59%). This is expected as the fully connected layer overfits to the training data if

Giải thích chi tiết:

1. Global Average Pooling (GAP) so với Fully Connected Layers:

- **Global Average Pooling (GAP)** và lớp fully connected đều thực hiện các biến đổi tuyến tính trên các bản đồ đặc trưng đã được vector hóa. Tuy nhiên, sự khác biệt giữa hai phương pháp này nằm ở cách chúng xử lý ma trận biến đổi. GAP sử dụng ma trận biến đổi đơn giản hơn, trong đó các giá trị không khác không chỉ xuất hiện trên các đường chéo chính. Trong khi đó, các lớp fully connected thường có các ma trận biến đổi đầy đặc hơn và cần được tối ưu hóa trong quá trình lan truyền ngược.

2. Thí nghiệm:

- Tác giả đã thực hiện các thí nghiệm để đánh giá hiệu quả của GAP so với lớp fully connected bằng cách thử nghiệm trên bộ dữ liệu CIFAR-10. Hai cấu hình mô hình được so sánh là:
 - Mô hình có lớp fully connected.
 - Mô hình sử dụng GAP thay cho lớp fully connected.
- Ngoài ra, cả hai mô hình này đều được kiểm tra với và không có dropout.

3. Kết quả:

- Bảng kết quả (Table 5) cho thấy rằng mô hình với fully connected có tỷ lệ lỗi cao hơn (11.59%) so với mô hình có dropout (10.88%). Tuy nhiên, mô hình sử dụng GAP đạt được kết quả tốt nhất với tỷ lệ lỗi thấp nhất (10.41%).
- Điều này cho thấy rằng GAP hoạt động hiệu quả hơn trong việc tránh quá khớp (overfitting), và là một regularizer tốt hơn so với việc chỉ sử dụng fully connected layers kết hợp với dropout.

no regularizer is applied. Adding dropout before the fully connected layer reduced the testing error (10.88%). Global average pooling has achieved the lowest testing error (10.41%) among the three.

We then explore whether the global average pooling has the same regularization effect for conventional CNNs. We instantiate a conventional CNN as described by Hinton et al. [5], which consists of three convolutional layers and one local connection layer. The local connection layer generates 16 feature maps which are fed to a fully connected layer with dropout. To make the comparison fair, we reduce the number of feature map of the local connection layer from 16 to 10, since only one feature map is allowed for each category in the global average pooling scheme. An equivalent network with global average pooling is then created by replacing the dropout + fully connected layer with global average pooling. The performances were tested on the CIFAR-10 dataset.

This CNN model with fully connected layer can only achieve the error rate of 17.56%. When dropout is added we achieve a similar performance (15.99%) as reported by Hinton et al. [5]. By replacing the fully connected layer with global average pooling in this model, we obtain the error rate of 16.46%, which is one percent improvement compared with the CNN without dropout. It again verifies the effectiveness of the global average pooling layer. Although it is slightly worse than the dropout regularizer result, we argue that it is not too demanding for linear convolution layers as it requires the linear layer to output the confidence maps of the categories.

Điều này chứng minh rằng GAP có tác dụng điều chuẩn (regularizer) hiệu quả, nhưng không vượt qua được dropout khi so sánh trực tiếp.
Tuy nhiên, tác giả cũng lưu ý rằng GAP có thể yêu cầu nhiều hơn đối với các lớp tích chập tuyến tính vì nó cần bộ lọc tuyến tính với hàm kích hoạt để mô hình hóa bản đồ tự tin (confidence maps) của các danh mục.

4.7 Visualization of NIN

We explicitly enforce feature maps in the last mlpconv layer of NIN to be confidence maps of the categories by means of global average pooling, which is possible only with stronger local receptive field modeling, e.g. mlpconv in NIN. To understand how much this purpose is accomplished, we extract and directly visualize the feature maps from the last mlpconv layer of the trained model for CIFAR-10.

Figure 4 shows some exemplar images and their corresponding feature maps for each of the ten categories selected from CIFAR-10 test set. It is expected that the largest activations are observed in the feature map corresponding to the ground truth category of the input image, which is explicitly enforced by global average pooling. Within the feature map of the ground truth category, it can be observed that the strongest activations appear roughly at the same region of the object in the original image. It is especially true for structured objects, such as the car in the second row of Figure 4. Note that the feature maps for the categories are trained with only category information. Better results are expected if bounding boxes of the objects are used for fine grained labels.

The visualization again demonstrates the effectiveness of NIN. It is achieved via a stronger local receptive field modeling using mlpconv layers. The global average pooling then enforces the learning of category level feature maps. Further exploration can be made towards general object detection. Detection results can be achieved based on the category level feature maps in the same flavor as in the scene labeling work of Farabet et al. [20].

5 Conclusions

We proposed a novel deep network called “Network In Network” (NIN) for classification tasks. This new structure consists of mlpconv layers which use multilayer perceptrons to convolve the input and a global average pooling layer as a replacement for the fully connected layers in conventional CNN. Mlpconv layers model the local patches better, and global average pooling acts as a structural regularizer that prevents overfitting globally. With these two components of NIN we demonstrated state-of-the-art performance on CIFAR-10, CIFAR-100 and SVHN datasets. Through visualization of the feature maps, we demonstrated that feature maps from the last mlpconv layer of NIN were confidence maps of the categories, and this motivates the possibility of performing object detection via NIN.

References

- [1] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

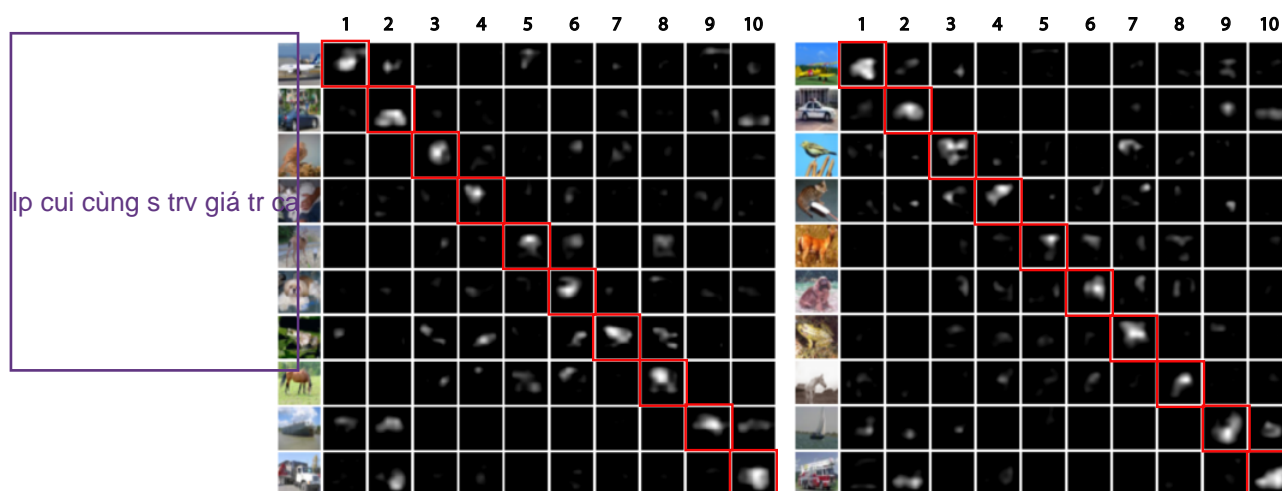


Figure 4: Visualization of the feature maps from the last mlpconv layer. Only top 10% activations in the feature maps are shown. The categories corresponding to the feature maps are: 1. airplane, 2. automobile, 3. bird, 4. cat, 5. deer, 6. dog, 7. frog, 8. horse, 9. ship, 10. truck. Feature maps corresponding to the ground truth of the input images are highlighted. The left panel and right panel are just different exemplars.

- [2] Y Bengio, A Courville, and P Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35:1798–1828, 2013.
- [3] Frank Rosenblatt. Principles of neurodynamics. perceptrons and the theory of brain mechanisms. Technical report, DTIC Document, 1961.
- [4] Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25*, pages 1106–1114, 2012.
- [5] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- [6] Quoc V Le, Alexandre Karpenko, Jiquan Ngiam, and Andrew Ng. Ica with reconstruction cost for efficient overcomplete feature learning. In *Advances in Neural Information Processing Systems*, pages 1017–1025, 2011.
- [7] Ian J Goodfellow. Piecewise linear multilayer perceptrons and dropout. *arXiv preprint arXiv:1301.5088*, 2013.
- [8] Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. *arXiv preprint arXiv:1302.4389*, 2013.
- [9] Çağlar Gülçehre and Yoshua Bengio. Knowledge matters: Importance of prior information for optimization. *arXiv preprint arXiv:1301.4083*, 2013.
- [10] Henry A Rowley, Shumeet Baluja, Takeo Kanade, et al. *Human face detection in visual scenes*. School of Computer Science, Carnegie Mellon University Pittsburgh, PA, 1995.
- [11] Matthew D Zeiler and Rob Fergus. Stochastic pooling for regularization of deep convolutional neural networks. *arXiv preprint arXiv:1301.3557*, 2013.
- [12] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- [13] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, volume 2011, 2011.
- [14] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. *arXiv preprint arXiv:1206.2944*, 2012.

- [15] Li Wan, Matthew Zeiler, Sixin Zhang, Yann L Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 1058–1066, 2013.
- [16] Mateusz Malinowski and Mario Fritz. Learnable pooling regions for image classification. *arXiv preprint arXiv:1301.3516*, 2013.
- [17] Nitish Srivastava and Ruslan Salakhutdinov. Discriminative transfer learning with tree-based priors. In *Advances in Neural Information Processing Systems*, pages 2094–2102, 2013.
- [18] Nitish Srivastava. *Improving neural networks with dropout*. PhD thesis, University of Toronto, 2013.
- [19] Ian J Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit number recognition from street view imagery using deep convolutional neural networks. *arXiv preprint arXiv:1312.6082*, 2013.
- [20] Clément Farabet, Camille Couprie, Laurent Najman, Yann Lecun, et al. Learning hierarchical features for scene labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:1915–1929, 2013.