

# Colinear Convolution Layer

Pang Liang

October 28, 2014

## 1 Cifar10 Database

The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images. (See <http://www.cs.toronto.edu/~kriz/cifar.html>)

The result of each approach propose by paper to reference [http://rodrigob.github.io/are\\_we\\_there\\_yet/build/classification\\_datasets\\_results.html](http://rodrigob.github.io/are_we_there_yet/build/classification_datasets_results.html).

Fig 1 are the classes in the dataset, as well as 10 random images from each.

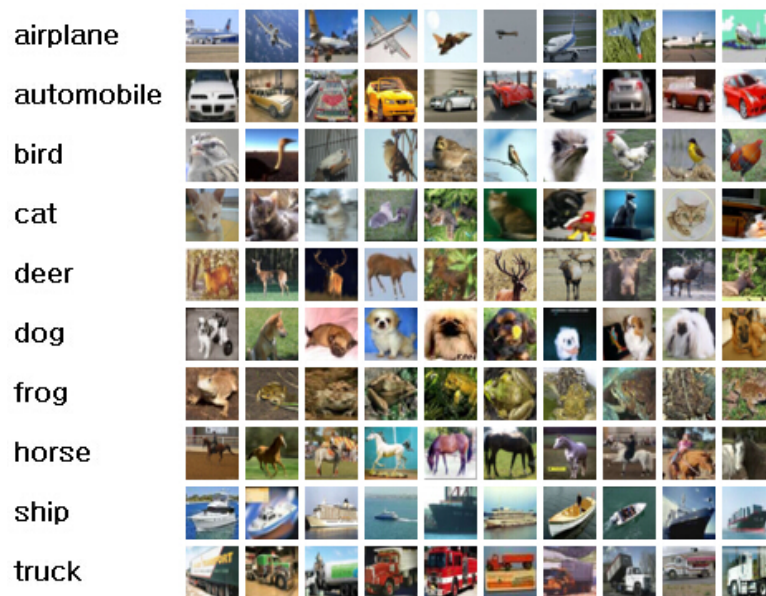


Figure 1: Example of cifar10 database.

## 2 Caffe Framework of Deep Learning

Caffe is a deep learning framework developed with cleanliness, readability, and speed in mind. (See <http://caffe.berkeleyvision.org>)

It implement many kind of layers list below:

- Vision Layers
  - Convolution ✓
  - Pooling ✓
  - Local Response Normalization (LRN) ✓
- Loss Layers
  - Softmax ✓
  - Sum-of-Squares / Euclidean
  - Hinge / Margin
  - Sigmoid Cross-Entropy
- Activation / Neuron Layers
  - ReLU / Rectified-Linear and Leaky-ReLU ✓
  - Sigmoid
  - TanH / Hyperbolic Tangent
- Common Layers
  - Inner Product ✓
  - Splitting
  - Flattening
  - Concatenation
  - Slicing

### 3 Fast Version of NN Structure for Cifar10

The Cifar10 model Fig 2 is a CNN that composes layers of convolution, pooling, rectified linear unit (ReLU) nonlinearities, and local contrast normalization with a linear classifier on top of it all.

The inputs of the model are not the raw images or raw images normalized into  $[0, 1]$ , but the images preprocess by the image mean operation. In other word, the pixel at same position average over all images, so the input pixel value range is  $[-127, 127]$ . In the end of the network we use Softmax loss as the classification loss.

## 4 Compare to Other Method

### 4.1 Network in Network

The goal of Network in Network is to nonlinear combine the different feature maps. It add two CCCP(Cascaded Cross Channel Parametric Pooling) Layers after Convolution Layer, and this kind of layer equivalent to the  $1 \times 1$  size Convolution Layer. Fig 3

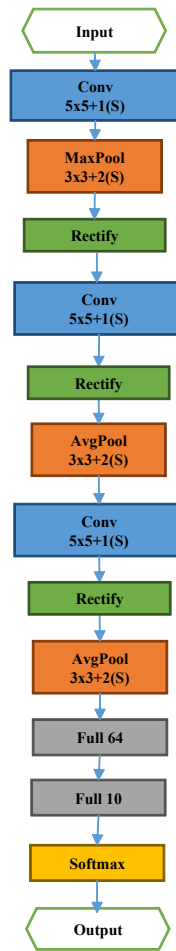


Figure 2: Cifar10 quick model.

I think the reason why add this kind of layer is to add nonlinear factor to the CNN output. Combining different feature maps using two layer neural network, is to say the linear CNN is not enough for the current application. So if we make CNN a little more complicate will get better result.

For more detail of the caffe configuration of NIN(CCCP Layer), see <https://gist.github.com/mavenlin/e56253735ef32c3c296d>.

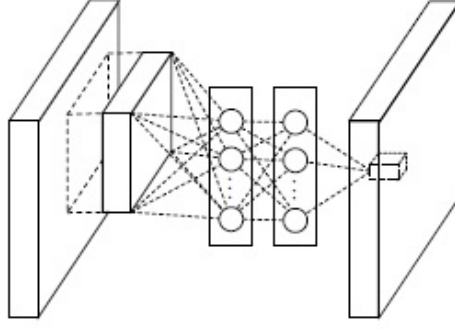


Figure 3: Network in network layer structure.

## 4.2 Drop Connection

## 4.3 GoogLeNet 2014

Google release it best ImageNet network structure in paper [??](#). It use  $1 \times 1$  size Convolution Layer too. But in this paper, they explain that this kind of Convolution Layers use to reduce the feature map number. This layer always puts in front of the  $5 \times 5$  or  $7 \times 7$  convolution layers, in order to reduce the input feature map of these large kernel layers. Fig 4

## 4.4 Spatially sparse CNN

Cifar-10 achieve 95.53% accuracy. See <http://www.kaggle.com/c/cifar-10/leaderboard/public>.

For more detail of his network structure is <http://www2.warwick.ac.uk/fac/sci/statistics/staff/academic-research/graham>.

# 5 Colinear Convolution Layer in Formula

## 5.1 Feed Forward

$$\begin{aligned} a_{in} &= e_i F_n e_i^T + S_n e_i^T + b_n \\ &= \sum_x \sum_y e_{ix} f_{xyn} e_{iy} + \sum_x s_{xn} e_{ix} + b_n \end{aligned} \quad (1)$$

## 5.2 Back Propagation

$$\sum_n \frac{\partial E}{\partial a_{in}} = \epsilon_{in} \quad (2)$$

- Error Propagation

$$\frac{\partial E}{\partial e_{ix}} = \sum_n \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial e_{ix}} = \sum_n \epsilon_{in} ((\sum_y (f_{xyn} + f_{yxn}) e_{iy}) + s_{xn}) \quad (3)$$



Figure 3: GoogLeNet network with all the bells and whistles

7

Figure 4: GoogLeNet structure.

- Weight Update

$$\frac{\partial E}{\partial f_{xyn}} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial f_{xyn}} = \sum_i \epsilon_{in} e_{ix} e_{iy} \quad (4)$$

$$\frac{\partial E}{\partial s_{xn}} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial s_{xn}} = \sum_i \epsilon_{in} e_{ix} \quad (5)$$

$$\frac{\partial E}{\partial b_n} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial b_n} = \sum_i \epsilon_{in} \quad (6)$$

See details in Fig 5.

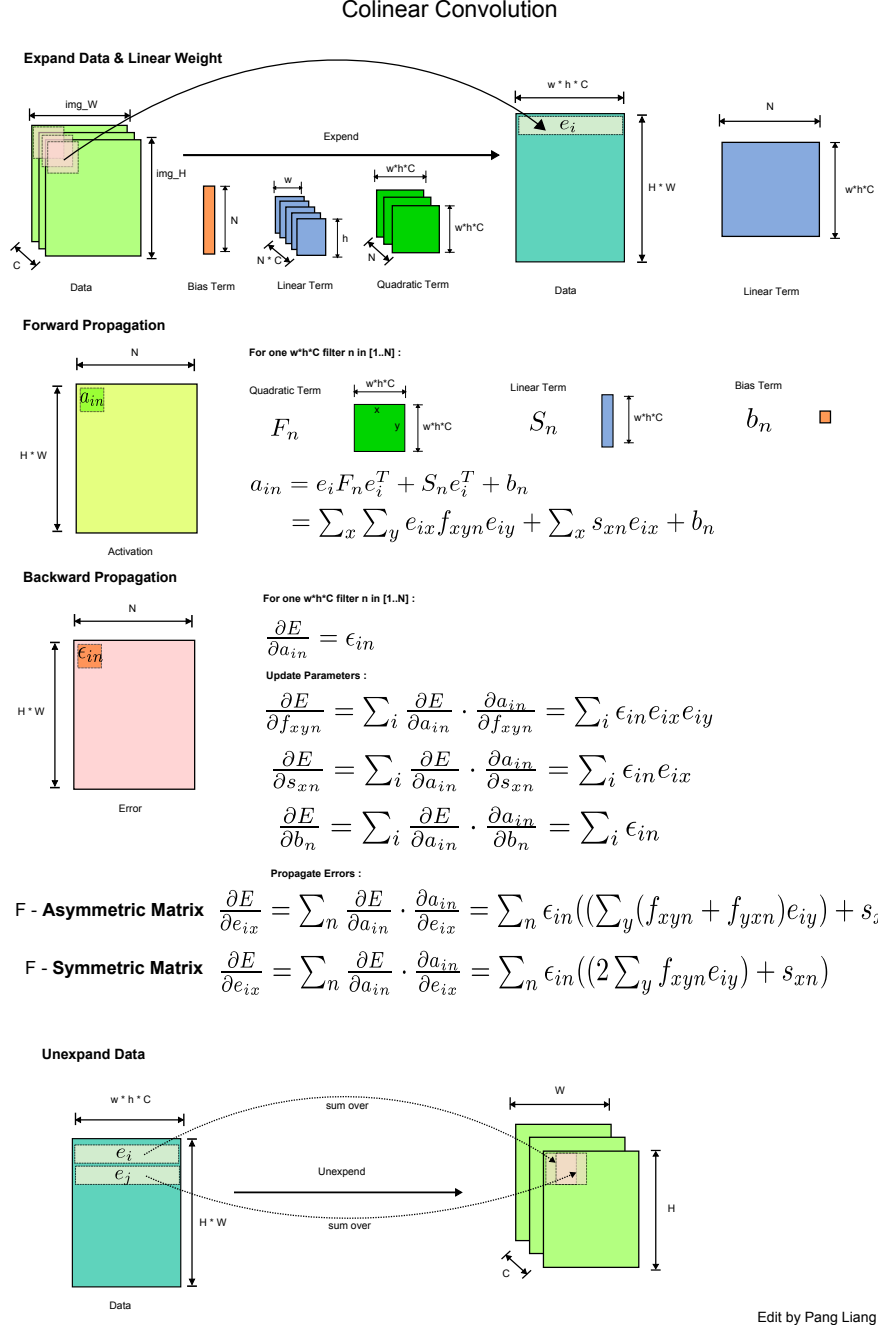


Figure 5: Colinear Convolution Layer.

## 6 Experiments

We replace the first Convolution layer to our Colinear Convolution layer. Weights in quadratic term are initialed in range of  $[0, 0.000001]$ , for the sake of the input range  $[-127, 127]$  and form of  $e * f * e$ . We want the activations of the layer not too large, so the weights in quadratic term are very small. For the same reason, the weight decay is set to 1000.

In Fig 6 and Fig 7 ,the number 16 or 32 denotes the first convolution layer output channel number. The lines entitled with Base mean that the result produce by Convolution layer, and others without produce by Colinear Convolution layer. The experiments show that the original configuration of the network is better than what we proposed. But there're some points we should optimize, which will show in next section.

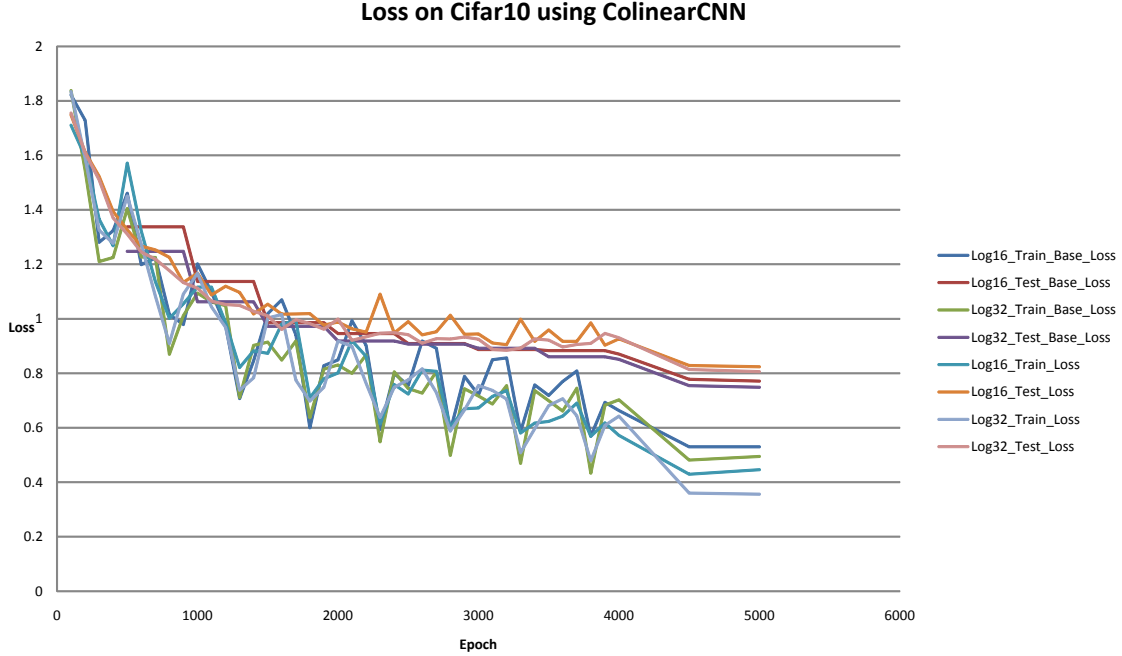


Figure 6: Train loss and Test loss of four net structures.

## 7 L1 Regularize on CCNN

Using L1 regularize on the parameters always produces a more sparse result, which meanings the bigger the regularizer  $\alpha$  the more 0 in parameters. Different from the L2 regularize  $\lambda||W||_2$ , L1 regularize has the form of  $\alpha||W||_1$ . L2 regularize can derived everywhere, see Eq 7.

$$\frac{\partial \lambda ||W||_2}{\partial W} = 2\lambda W \quad (7)$$

While L1 doesn't has this good property. Its derivative is Eq 8.

$$\frac{\partial \alpha ||W||_1}{\partial W} = \begin{cases} -\alpha & W < 0 \\ undefine & W = 0 \\ \alpha & W > 0 \end{cases} \quad (8)$$

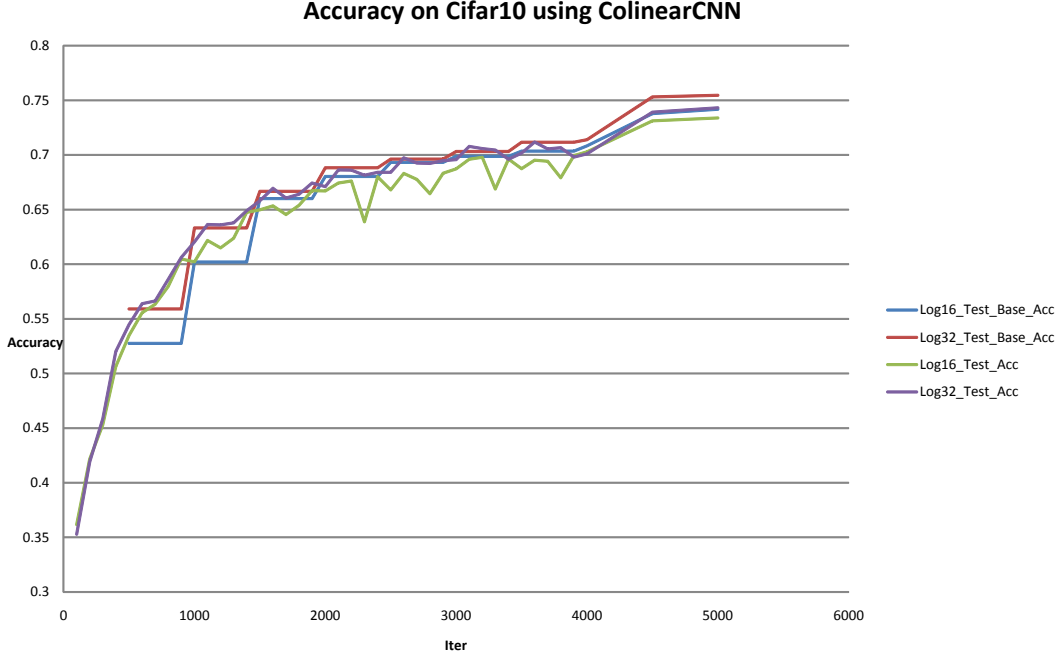


Figure 7: Test Accuracy of four net structures.

So in order to implement L1 regularize, we need to know  $W$ 's sign. For instance one parameter  $w_t$  at  $t$ 's iterator, then we want to calculate next update of  $w_t$  denoted as  $w_{t+1}$ .

$$w_{t+1} = w_t + \Delta w - \text{sign}(w)\alpha \quad (9)$$

$$w_{t+1}^{\hat{}} = w_t + \Delta w \quad (10)$$

then we need  $w_{t+1}$  and  $w_t + \Delta w$  have the same sign. So the last term  $-\text{sign}(w)\alpha$  will not change the sign. So we use Eq 11 for implement.

$$w_{t+1} = \begin{cases} \max(w_{t+1}^{\hat{}} - \alpha, 0) & w_{t+1}^{\hat{}} \geq 0 \\ \min(w_{t+1}^{\hat{}} + \alpha, 0) & w_{t+1}^{\hat{}} < 0 \end{cases} \quad (11)$$

## 8 Symmetric Strategy on CCNN

Changing Asymmetric Matrix  $F$  to Symmetric Matrix is natural to us, because the relation between  $e_{ix}$  and  $e_{iy}$  should be same. That means  $f_{xyn} = f_{yxn}$ . So we decrease parameter number from  $N^2$  to  $N * (N + 1)/2$ .

Then we need modify some formula in BP process.

- Error Propagation

$$\frac{\partial E}{\partial e_{ix}} = \sum_n \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial e_{ix}} = \sum_n \epsilon_{in} \left( \left( \sum_y 2 \cdot \begin{cases} f_{xyn} & x \leq y \\ f_{yxn} & x > y \end{cases} \cdot e_{iy} \right) + s_{xn} \right) \quad (12)$$



- Weight Update

$$\frac{\partial E}{\partial f_{xyn}} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial f_{xyn}} = \begin{cases} 2 \sum_i \epsilon_{in} e_{ix} e_{iy} & x < y \\ \sum_i \epsilon_{in} e_{ix} e_{iy} & x = y \end{cases} \quad (13)$$

## 9 Factorize Quadratic Matrix

Not only assume quadratic matrix is symmetric, but also low rank. So we can factorize this matrix to two low rank matrix. That is to say:

$$F = HH^T. \quad (14)$$

where matrix  $H$  is a  $N * K$  matrix. Then our parameters number decrease to  $N * K$ .

So the next steps is to modify FP and BP procession.

- Forward Propagation

$$\begin{aligned} a_{in} &= e_i H_n H_n^T e_i^T + S_n e_i^T + b_n \\ &= \sum_x \sum_y \sum_k e_{ix} h_{xkn} h_{ykn} e_{iy} + \sum_x s_{xn} e_{ix} + b_n \\ &= \sum_k \left( \sum_x e_{ix} h_{xkn} \right)^2 + \sum_x s_{xn} e_{ix} + b_n \end{aligned} \quad (15)$$

- Error Propagation

$$\begin{aligned} \frac{\partial E}{\partial e_{ix}} &= \sum_n \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial e_{ix}} \\ &= \sum_n \epsilon_{in} \left( \left( \sum_y \sum_k 2 h_{xkn} h_{ykn} e_{iy} \right) + s_{xn} \right) \\ &= 2 \sum_n \epsilon_{in} \sum_y e_{iy} \sum_k h_{xkn} h_{ykn} + \sum_n \epsilon_{in} s_{xn} \end{aligned} \quad (16)$$

- Weight Update

$$\begin{aligned} \frac{\partial E}{\partial h_{xkn}} &= \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial h_{xkn}} \\ &= \sum_i \epsilon_{in} 2 \sum_y e_{ix} e_{iy} h_{ykn} \\ &= 2 \sum_i \epsilon_{in} e_{ix} \sum_y e_{iy} h_{ykn} \end{aligned} \quad (17)$$

$$\frac{\partial E}{\partial s_{xn}} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial s_{xn}} = \sum_i \epsilon_{in} e_{ix} \quad (18)$$

$$\frac{\partial E}{\partial b_n} = \sum_i \frac{\partial E}{\partial a_{in}} \cdot \frac{\partial a_{in}}{\partial b_n} = \sum_i \epsilon_{in} \quad (19)$$

Consider a special condition that set  $K = 1$ . That means our matrix can be factorize into two vectors dot product. In this situation we can use a special active function and convolutional layer to mimic that. It's easy to prove that this kind of factorization is equivalent to the convolutional layer add a quadratic active function -  $x^2$ .

$$\sum_k (\sum_x e_{ix} h_{xkn})^2 = \sum_x \sum_y \sum_k e_{ix} h_{xkn} h_{ykn} e_{iy}. \quad (20)$$

## 10 Mask the Parameters

Maybe this idea can be compared with dropout connections.

## 11 Restriction and Regularization of CCNN

1. Change the Asymmetric to Symmetric, in order to decrease the parameter numbers.
2. Change the L2 regularize to L1 regularize, for the 2 order relation is weaker than the 1 order relation and we need not that much parameters to model it.
3. Using mask matrix to erase the relations cross the channels.