

## Lab 3: CNN on Cifar-10

### Lab Objective:

In this lab, you will be asked to build popular network architecture (*Network In Network*, NIN) [1], and train it on Cifar-10 dataset. Moreover, you need to use data augmentation and Dropout [2] during training.

### Turn in:

1. Experiment Report (7/24(一) 12:00)
2. Demo date (7/24(一) 下課後)

### Requirements:

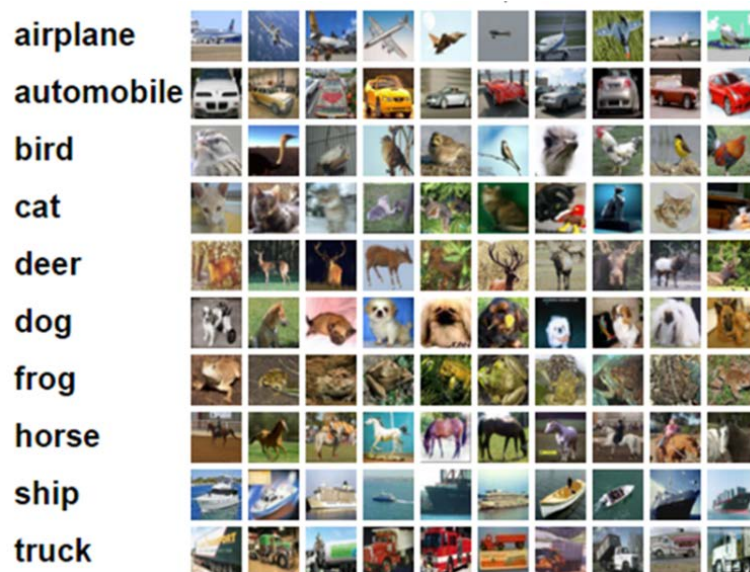
- Implement “**Network In Network**” (NIN) [1] convolutional architecture
- Implement data augmentation: translation and horizontal flipping
- Use “**Dropout**” [2] in NIN
- Train NIN+Dropout with/without data augmentation

### Environment:

- Cifar-10 dataset

The CIFAR-10 dataset consists of 60000  $32 \times 32$  color images (RGB) in **10** classes, with 6000 images per class. There are 50000 training images and 10000 test images.

Download: <https://www.cs.toronto.edu/~kriz/cifar.html>



## Sample Code

There are many cifar-10 sample codes for [tensorflow](#):

*git clone <https://github.com/tensorflow/models.git>*

There are many cifar-10 sample codes for [pytorch](#):

*git clone <https://github.com/kuangliu/pytorch-cifar>*

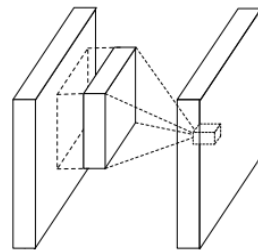
## Lab Description:

### ■ Network In Network (NIN)

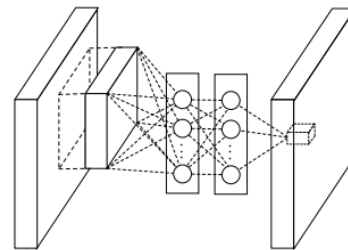
#### ■ Enhance model discriminability for local patches within the receptive field

Traditional CNN: 3x3 conv + ReLU

NIN: 3x3 conv + ReLU + 1x1 conv + ReLU + 1x1 conv + ReLU



(a) Linear convolution layer



(b) Mlpconv layer

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

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

### ■ Full NIN architecture used in Cifar-10

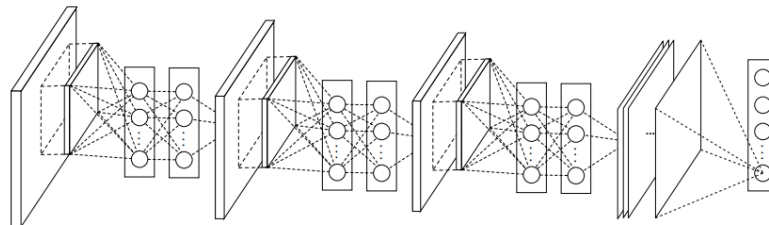


Figure 2: The overall structure of Network In Network. In this paper the NINs include the stacking of three mlpconv layers and one global average pooling layer.

### ■ Architecture Details:

Conv1	filter size = 5x5, # of filter =192, pad = 2, stride = 1	Act.=ReLU
mlp 1	filter size = 1x1, # of filter =160, pad = 0, stride = 1	Act.=ReLU
mlp 2	filter size = 1x1, # of filter =96, pad = 0, stride = 1	Act.=ReLU
Pool 1	3x3 max pooling, stride = 2, pad = 1 ('same')	
	Dropout 0.5	
Conv2	filter size = 5x5, # of filter =192, pad = 2, stride = 1	Act.=ReLU
mlp 2-1	filter size = 1x1, # of filter =192, pad = 0, stride = 1	Act.=ReLU
mlp 2-2	filter size = 1x1, # of filter =192, pad = 0, stride = 1	Act.=ReLU
Pool 2	3x3 max (avg) pooling, stride = 2, pad = 1 ('same')	
	Dropout 0.5	
Conv3	filter size = 3x3, # of filter =192, pad = 1, stride = 1	Act.=ReLU
mlp 3-1	filter size = 1x1, # of filter =192, pad = 0, stride = 1	Act.=ReLU
mlp 3-2	filter size = 1x1, # of filter =10, pad = 0, stride = 1	Act.=ReLU
Global Pool	8x8 average pooling, stride =1, pad = 0 ('same')	
	Softmax	

### ■ Data augmentation: Translation and Horizontal flipping:



Original



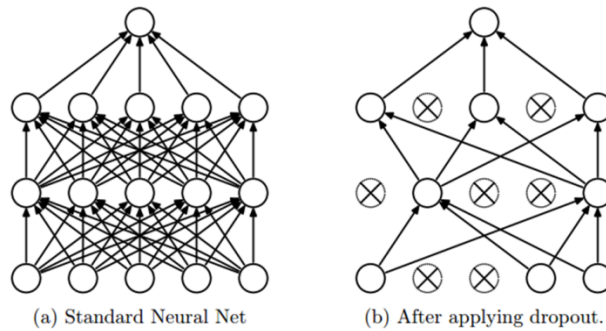
Translation



Horizontal flipping

### ■ Dropout

Dropout is a technique for addressing this problem. The key idea is to **randomly drop units** (along with their connections) from the neural network during training. This prevents units from **co-adapting** too much. During training, dropout samples from an exponential number of different “thinned” networks. At test time, it is easy to approximate the effect of averaging the predictions of all these thinned networks by simply using a single unthinned network that has **smaller weights**. This significantly **reduces overfitting** and gives major improvements over other regularization methods.



## ■ Network with Dropout

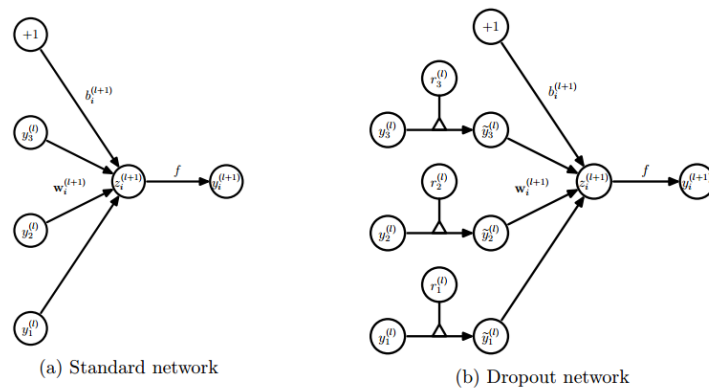
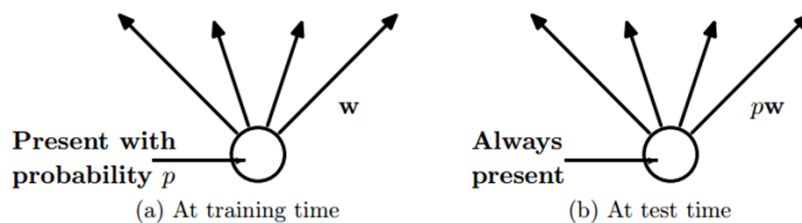


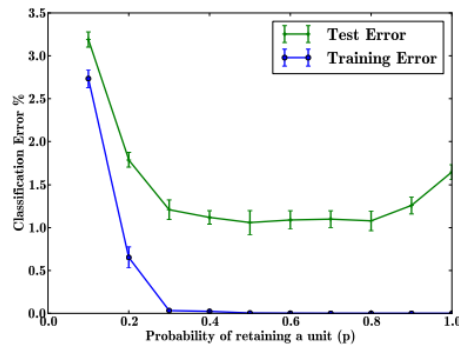
Figure 3: Comparison of the basic operations of a standard and dropout network.

## ■ Dropout at testing



## ■ The effect of Dropout rate

If the architecture is held constant, having a small  $p$  means very few units will turn on during training. It can be seen that this has led to *underfitting* since the training error is also high. We see that as  $p$  increases, the error goes down. It becomes flat when  $0.4 \leq p \leq 0.8$  and then increases as  $p$  becomes close to 1.



(a) Keeping  $n$  fixed.

### Implementation Details:

- Training Hyperparameters:
  - Method: SGD with momentum
  - Mini-batch size: **128** (391 iterations for each epoch)
  - Total epochs: **164**, momentum **0.9** (if you use momentum SGD)
  - Initial learning rate: **0.1**, divide by 10 at 81, 122 epoch
  - Loss function: cross-entropy
- Data augmentation parameters:
  - Translation: Pad **4** zeros in each side and random cropping back to 32x32 size
  - Horizontal flipping: With probability **0.5**

### Methodology:

- 91.10% accuracy with data augmentation in my implementation
- 89.23% accuracy without data augmentation

### Extra Bonus:

- All convolutions NIN (remove pooling layers)
- Reproduce the experiment of Dropout rate in [2].

### References:

- [1] Lin, M., Chen, Q., & Yan, S. (2013). Network in network. *arXiv preprint arXiv:1312.4400*.
- [2] Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1), 1929-1958.
- [3] [https://www.tensorflow.org/tutorials/deep\\_cnn/](https://www.tensorflow.org/tutorials/deep_cnn/)

Report Spec: [black: Demo, Gray: No Demo]

1. Introduction (15%)

2. Experiment setup (15%)

- The detail of your model
- Report all your training hyper-parameters

3. Result

- The comparison between with and without data augmentation
  - Final Test error (10%, 15%)
  - Training loss curve (you need to record training loss every epoch) (10%, 15%)
  - Test error curve (you need to record test error every epoch) (10%, 15%)

4. Discussion (20%, 25%)

Demo (20%) [抽 13 人 DEMO]

-----實驗結果標準 (with data augmentation)-----

Accuracy: (92.0~90.0)% = 100%

Accuracy: (90.0~87.0)% = 90%

Accuracy: (87.0~84.0)% = 80%

Accuracy: below 84.0% = 70%

Accuracy: 10% = 0%

評分標準: 40%\*實驗結果 + 60%\*(報告+DEMO)

## Data loading

<https://www.cs.toronto.edu/~kriz/cifar.html>

### Python / Matlab versions

I will describe the layout of the Python version of the data

The archive contains the files `data_batch_1`, `data_batch_2`, ..., `data_batch_10`. Each file returns a dictionary:

```
def unpickle(file):
    import cPickle
    fo = open(file, 'rb')
    dict = cPickle.load(fo)
    fo.close()
    return dict
```

Loaded in this way, each of the batch files contains a dictionary with two keys:

- **data** -- a 10000x3072 `numpy` array of `uint8s`. Each image is stored in row-major order, so that the first 10000 rows are the first 10000 images.
- **labels** -- a list of 10000 numbers in the range 0-9.

```
10 def unpickle(file):
11     import cPickle
12     fo = open(file, 'rb')
13     dict = cPickle.load(fo)
14     fo.close()
15     if 'data' in dict:
16         dict['data'] = dict['data'].reshape((-1, 3, 32, 32)).swapaxes(1, 3).swapaxes(1, 2).reshape((-1, 32*32*3))
17     return dict
18
19
20 def load_data_one(f):
21     batch = unpickle(f)
22     data = batch['data']
23     labels = batch['labels']
24     print("Loading %s: %d" % (f, len(data)))
25     return data, labels
26
27 def load_data(files, data_dir, label_count):
28     data, labels = load_data_one(data_dir + '/' + files[0])
29     for f in files[1:]:
30         data_n, labels_n = load_data_one(data_dir + '/' + f)
31         data = np.append(data, data_n, axis=0)
32         labels = np.append(labels, labels_n, axis=0)
33     labels = np.array([ [ float(i == label) for i in xrange(label_count) ] for label in labels ])
34     return data, labels
```

```
94 ###
95 data_dir = './cifar-10-batches-py'
96 image_size = 32
97 image_dim = image_size * image_size * 3
98 meta = unpickle(data_dir + '/batches.meta')
99 label_names = meta['label_names']
100 label_count = len(label_names)
101
102 train_files = [ 'data_batch_%d' % d for d in xrange(1, 6) ]
103 train_data, train_labels = load_data(train_files, data_dir, label_count)
104 #print("Train:", np.shape(train_data), np.shape(train_labels))
105 #print("Test:", np.shape(test_data), np.shape(test_labels))
106 test_data, test_labels = load_data([ 'test_batch' ], data_dir, label_count)
107 pi = np.random.permutation(len(train_data))
108 train_data, train_labels = train_data[pi], train_labels[pi]
109 train_data = train_data.reshape((-1, 32, 32, 3))
```

### data augmentation

[https://github.com/JiaRenChang/DLcourse\\_NCTU/blob/master/data\\_aug.py](https://github.com/JiaRenChang/DLcourse_NCTU/blob/master/data_aug.py)

Usage:

Import data\_aug

```
cropped_data = data._random_crop(xs_, [32, 32, 3], padding=4)
```

```
flipped_data = data._random_flip_leftright(cropped_data)
```

### Save session in tensorflow

[https://www.tensorflow.org/programmers\\_guide/variables](https://www.tensorflow.org/programmers_guide/variables)

```
saver = tf.train.Saver()
```

```
save_path = saver.save(session, 'NIN_%d.ckpt' % epoch)
```

### Save in pytorch

```
savefilename = 'yourfilename.tar'
```

```
torch.save({'state_dict': model.state_dict()}, savefilename)
```



## Tricks for LAB 1

### 1. Data preprocessing

#### ■ Color normalization

Normalize each color channel (compute from entire CIFAR10 training set)

$$\text{Mean} \quad \begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{pmatrix} 125.3 \\ 123.0 \\ 113.9 \end{pmatrix}$$

$$\text{Standard deviation} \quad \begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{pmatrix} 63.0 \\ 62.1 \\ 66.7 \end{pmatrix}$$

This trick provides about 4~5% improvement

Without data preprocessing:

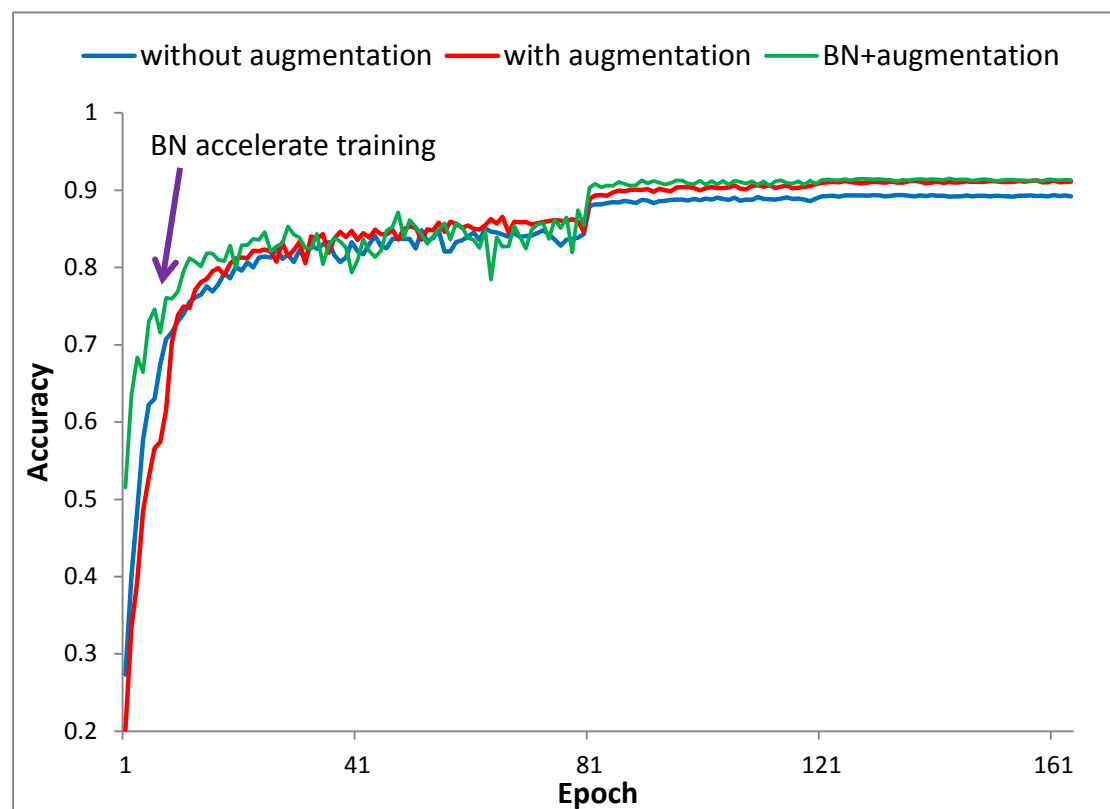
Without data augmentation: ~84%

With data augmentation: ~85%

With data preprocessing:

Without data augmentation: 89.23%

With data augmentation: 91.10%



## 2. Learning rate schedule

0.1: 1~80    0.01: 81~160    0.001: 161~200

Tensorflow: Let learning rate is placeholder:

```
learning_rate = tf.placeholder(tf.float32)
```

```
lr = 0.1
```

```
if epoch >= 81: lr = 0.01
```

```
if epoch >= 122: lr = 0.001
```

```
batch_res = session.run([ train_step, accuracy, cross_entropy],  
                        feed_dict = {x: xs_, y_: ys_, learning_rate: lr, keep_prob: 0.5 })
```

-----

PyTorch:

```
def adjust_learning_rate(optimizer, epoch):  
    if epoch <= 80:  
        lr = 0.1  
    elif epoch <= 122:  
        lr = 0.01  
    else:  
        lr = 0.001  
    for param_group in optimizer.param_groups:  
        param_group['lr'] = lr
```

### 3. Weight Decay

Weight decay = 0.0001

Tensorflow:

```
l2 = tf.add_n([tf.nn.l2_loss(var) for var in tf.trainable_variables()])  
# optimizer SGD  
train_step = tf.train.MomentumOptimizer(learning_rate, 0.9,  
use_nesterov=True).minimize(cross_entropy+ l2 * Weight decay)
```

Pytorch:

```
ptimizer = torch.optim.SGD(model.parameters(), args.lr,  
                             momentum=True,  
                             weight_decay=0.0001)
```

### 4. Weight initialization in LAB 1

First conv layer: Random\_normal(stddev=0.01)

Others: Random\_normal(stddev=0.05)

### 5. Use Nesterov momentum

Tensorflow:

```
tf.train.MomentumOptimizer(learning_rate, 0.9, use_nesterov=True)
```

Pytorch:

```
ptimizer = torch.optim.SGD(model.parameters(), args.lr,  
                             momentum=True,  
                             weight_decay=0.0001)
```