

SHARIF UNIVERSITY OF TECHNOLOGY



COMPUTER VISION

---

## HW05: Scene Recognition Using Neural Networks

---

*Author:*  
Saeed Razavi (98106542)

July 16, 2022

## question1

in this notebook we want to train cnn models with different architectures .

- PART 1 : using **first convolutional layer and two last fully connected layer of Alexnet architecture** . note that we use **max-pooling** layer with **kernel-size=4** and **stride=4** in first layer .
- PART 2 : using **three first convolutions layer and three last fully connected layer of Alexnet architecture** . in the **third layer** , we change the conv2d layer with depth of 384 to 256 .
- PART 3 : using **whole Alexnet architecture** and just change the last fully-connected layer with depth of 1000 to 15 but **without any pre-trained weights**. more specifically , we don't use pre-trained parameters that derived for architecture from imagenet dataset .
- PART 4 : using **whole Alexnet architecture** and just change the last fully-connected layer with depth of 1000 to 15 . we **use pre-trained weights** for all the layers expect for new classifier layer .during training , **we just let the parameters of new classifier layer to update** in each epoch and freeze all other parameters.
- PART 5 : using **whole Alexnet architecture** and just change the last fully-connected layer with depth of 1000 to 15 . we **use pre-trained weights** for all the layers expect for new classifier layer .during training , **we let all parameters of all layers to update in each epoch** .

## hyperparameters

we determine hyperparameters as follows :

```
1 config = {
2     '''
3     note that all files are saved in folders.
4     Pictures of the same category are stored in each folder.
5     The folder name is class name
6
7     '''
8
9
10    'extracted_root' : "/content/Data", #path of extracted zip data
11    'train_data_path' : "/content/Data/Train" , #path of extracted train data
12    'test_data_path' : "/content/Data/Test" , # path of extracted test data
13    'batch_size' : 64 ,
14    'num_workers' : 1 ,
15    'num_classes' : 15 ,
16    'epoches' : 50 ,
17    'learning_rate' : 0.0005
18    'optimizer' : "torch.optim.SGD"
19 }
```

Listing 1: hyperparameters

- **batch-size** : set to **64** because for greater value the less update we have for the weights (doesn't converge well) and for smaller value we lost the generality of the model (the gradient is not trust-able)
- **learning rate** : set to **0.0005** and gradually decrease in next epoch
- **num-classes** : set to **15** because our classification problem has 15 class
- **epoches** : number of epoches set to **50** for this problem
- **num-workers** : indicate how multi-processing we have in each batch . is set it to **1**
- **train\_data** – path : *path of extracted train data that we have pass to* **datasets.ImageFolder**
- **test-data-path** : path of extracted test data that we have pass to **torchvision.datasets.ImageFolder**

# Transforms

we set different transforms for training and test set

- training set : for this , we define 4 different transform in a sequential structure .
  1. resize each image (**transforms.Resize**) to  $227 \times 227$
  2. augment training set using **transforms.RandomHorizontalFlip** to mirror each image horizontally . note that because our objective is **scene recognition** , many of transforms like rotating , mirror vertically and ... lead to bad training and bad dataset . so we only use **transforms.RandomVerticalFlip**. we set the p(probability) to 1 to have one new training dataset alongside the original training dataset
  3. convert dataset to tensor
  4. normalize dataset based on alexnet mean and std
- test set : all of test-set transforms are like training-set transforms except that we don't need augmentation in test-set

```
1 #define normalization transform
2 normalize = transforms.Normalize(
3     mean=[0.4914, 0.4822, 0.4465],
4     std=[0.2023, 0.1994, 0.2010],)
5
6
7 #define train_set transform
8 train_transform1 = transforms.Compose([
9     transforms.Resize((227,227)),
10    transforms.ToTensor(),
11    normalize,])
12
13 #-----
14
15 train_transform2 = transforms.Compose([
16     transforms.Resize((227,227)),
17     transforms.RandomHorizontalFlip(p=1),
18     transforms.ToTensor(),
19     normalize,])
20 # #-----
21
22 #define test_set transform
23 test_transform = transforms.Compose([
24     transforms.Resize((227,227)),
25     transforms.ToTensor(),
26     normalize,])
```

Listing 2: Transforms

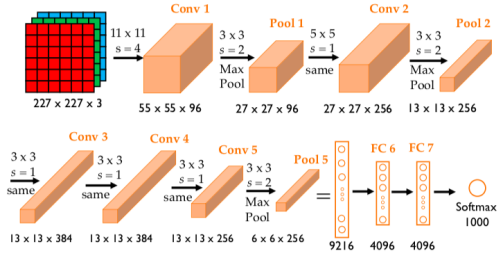
## splitting data and augmentation

one good way to split data and pass it to **dataloader** is to use **torchvision.datasets.ImageFolder** . this function get images with their label without using any label file . **Imagefolder** assumes that all files are saved in folders ,Pictures of the same category are stored in each folder and The folder name is class name . we also use **torch.utils.data.ConcatDataset** for concatenating the new training dataset created by **transforms.RandomHorizontalFlip(p=1)** with original training dataset

```
1 train_data_aug1 = torchvision.datasets.ImageFolder(root=config['train_data_path'], transform=
2     train_transform1)
3 all_datasets.append(train_data_aug1)
4 train_data_aug2 = torchvision.datasets.ImageFolder(root=config['train_data_path'], transform=
5     train_transform2)
6 all_datasets.append(train_data_aug2)
7
8 train_data_total = torch.utils.data.ConcatDataset(all_datasets)
```

Listing 3: splitting data and augmentation

## alexnet architecture



(a) frame450

	Layer	Feature Map	Size	Kernel Size	Stride	Activation
	Input	Image	1	227x227x3	-	-
1	Convolution	96	55 x 55 x 96	11x11	4	relu
	Max Pooling	96	27 x 27 x 96	3x3	2	relu
2	Convolution	256	27 x 27 x 256	5x5	1	relu
	Max Pooling	256	13 x 13 x 256	3x3	2	relu
3	Convolution	384	13 x 13 x 384	3x3	1	relu
4	Convolution	384	13 x 13 x 384	3x3	1	relu
5	Convolution	256	13 x 13 x 256	3x3	1	relu
	Max Pooling	256	6 x 6 x 256	3x3	2	relu
6	FC	-	9216	-	-	relu
7	FC	-	4096	-	-	relu
8	FC	-	4096	-	-	relu
Output	FC	-	1000	-	-	Softmax

(b) frame270

Figure 1

## part1

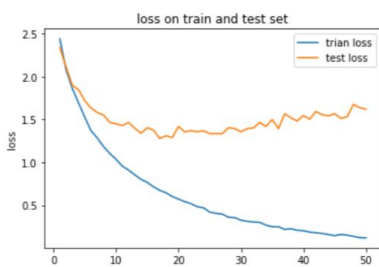
using first convolutional layer and two last fully connected layer of Alexnet architecture . note that we use **max-pooling** layer with **kernel-size=4** and **stride=4** in first layer . we also add new layer **LocalResponseNorm(2)** .this operation helps us to reduce input image noises below you can see the structure of network :

```

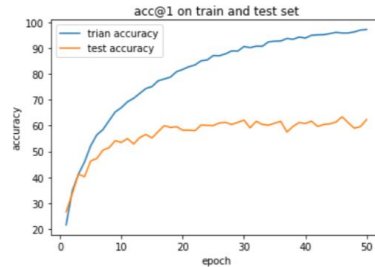
1 class AlexNet(nn.Module):
2     def __init__(self, num_classes=15):
3         super(AlexNet, self).__init__()
4         self.layer1 = nn.Sequential(
5             nn.Conv2d(3, 96, kernel_size=11, stride=4, padding=0),
6             # nn.BatchNorm2d(96),
7             nn.LocalResponseNorm(2), #adding new additional layer
8             nn.ReLU(),
9             nn.MaxPool2d(kernel_size = 4, stride = 4))
10        self.fc1 = nn.Sequential(
11            nn.Dropout(0.5),
12            nn.Linear(16224, 4096),
13            nn.ReLU())
14        self.fc2= nn.Sequential(
15            nn.Linear(4096, num_classes))
16
17        def forward(self, x):
18            out = self.layer1(x)
19            out = out.reshape(out.size(0), -1)
20            out = self.fc1(out)
21            out = self.fc2(out)
22            return out

```

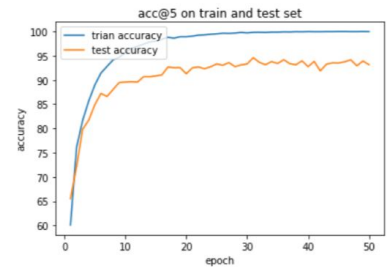
Listing 4: structure of part1



(a) loss



(b) acc@1



(c) acc@5

Figure 2

*Results :*

*maximum acc@1 on test set converges to : 63.4*

*maximum acc@5 on test set converges to : 93*

*maximum acc@1 on train set converges to : 100*

*maximum acc@5 on train set converges to : 97.2*

conclusions :

as we can see there is no overfitting in training model . there is a gap between test loss and train loss after some epoches as we expected. also the curve of acc@1-vs-epoch Meet our expectations because there is a gap between test's acc@1 and train's acc@1 over epoches.

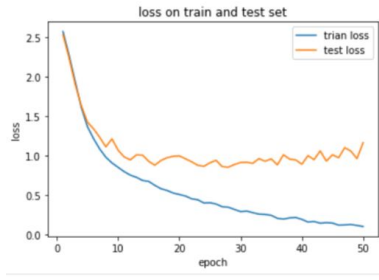
## part2

using three first convolutions layer and three last fully connected layer of Alexnet architecture . in the **third layer** , we change the conv2d layer with depth of 384 to 256 .

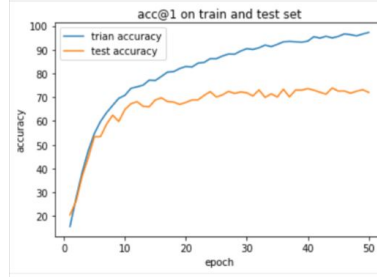
we also add new layer **LocalResponseNorm(2)** .this operation helps us to reduce input image noises below you can see the structure of network :

```
1 class AlexNet(nn.Module):
2     def __init__(self, num_classes=15):
3         super(AlexNet, self).__init__()
4         self.layer1 = nn.Sequential(
5             nn.Conv2d(3, 96, kernel_size=11, stride=4, padding=0),
6             # nn.BatchNorm2d(96),
7             nn.LocalResponseNorm(2), #adding new additional layer
8             nn.ReLU(),
9             nn.MaxPool2d(kernel_size = 3, stride = 2))
10        self.layer2 = nn.Sequential(
11            nn.Conv2d(96, 256, kernel_size=5, stride=1, padding=2),
12            nn.BatchNorm2d(256),
13            nn.ReLU(),
14            nn.MaxPool2d(kernel_size = 3, stride = 2))
15        self.layer3 = nn.Sequential(
16            nn.Conv2d(256, 256, kernel_size=3, stride=1, padding=1),
17            nn.BatchNorm2d(256),
18            nn.ReLU(),
19            nn.MaxPool2d(kernel_size = 2, stride = 2))
20        self.fc = nn.Sequential(
21            nn.Dropout(0.5),
22            nn.Linear(9216, 4096),
23            nn.ReLU())
24        self.fc1 = nn.Sequential(
25            nn.Dropout(0.5),
26            nn.Linear(4096, 4096),
27            nn.ReLU())
28        self.fc2= nn.Sequential(
29            nn.Linear(4096, num_classes))
30
31        def forward(self, x):
32            out = self.layer1(x)
33            out = self.layer2(out)
34            out = self.layer3(out)
35            out = out.reshape(out.size(0), -1)
36            out = self.fc(out)
37            out = self.fc1(out)
38            out = self.fc2(out)
39            return out
40
```

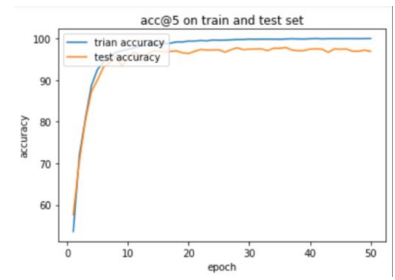
Listing 5: structure of part2



(a) loss



(b) acc@1



(c) acc@5

Figure 3

#### Results :

maximum acc@1 on test set converges to : 73.91

maximum acc@5 on test set converges to : 93.5

maximum acc@1 on train set converges to : 97.27

maximum acc@5 on train set converges to : 100

---

#### conclusions :

as we can see there is no overfitting in training model , in addition the curve of loss-vs-epoch Meet our expectations because at first , loss of test is bigger than train and over epoches , the train loss becomes less and less but test loss becomes approximately constant after the 20th epoch . there is also a gap between test loss and train loss after some epoches. also the curve of acc@1-vs-epoch Meet our expectations because there is a gap between test's acc@1 and train's acc@1 over epoches.

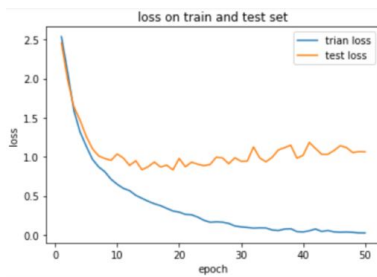
## part3

whole alexnet architecture without using trained weights

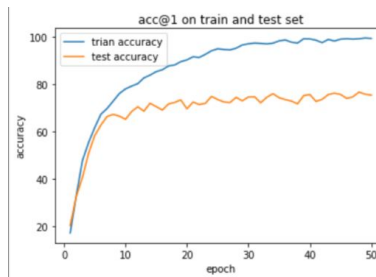
below you can see the structure of network :

```
1 class AlexNet(nn.Module):
2     def __init__(self, num_classes=15):
3         super(AlexNet, self).__init__()
4         self.layer1 = nn.Sequential(
5             nn.Conv2d(3, 96, kernel_size=11, stride=4, padding=0),
6             nn.BatchNorm2d(96),
7             nn.ReLU(),
8             nn.MaxPool2d(kernel_size = 3, stride = 2))
9         self.layer2 = nn.Sequential(
10            nn.Conv2d(96, 256, kernel_size=5, stride=1, padding=2),
11            nn.BatchNorm2d(256),
12            nn.ReLU(),
13            nn.MaxPool2d(kernel_size = 3, stride = 2))
14        self.layer3 = nn.Sequential(
15            nn.Conv2d(256, 384, kernel_size=3, stride=1, padding=1),
16            nn.BatchNorm2d(384),
17            nn.ReLU())
18        self.layer4 = nn.Sequential(
19            nn.Conv2d(384, 384, kernel_size=3, stride=1, padding=1),
20            nn.BatchNorm2d(384),
21            nn.ReLU())
22        self.layer5 = nn.Sequential(
23            nn.Conv2d(384, 256, kernel_size=3, stride=1, padding=1),
24            nn.BatchNorm2d(256),
25            nn.ReLU(),
26            nn.MaxPool2d(kernel_size = 3, stride = 2))
27        self.fc = nn.Sequential(
28            nn.Dropout(0.5),
29            nn.Linear(9216, 4096),
30            nn.ReLU())
31        self.fc1 = nn.Sequential(
32            nn.Dropout(0.5),
33            nn.Linear(4096, 4096),
34            nn.ReLU())
35        self.fc2 = nn.Sequential(
36            nn.Linear(4096, num_classes))
37
38    def forward(self, x):
39        out = self.layer1(x)
40        out = self.layer2(out)
41        out = self.layer3(out)
42        out = self.layer4(out)
43        out = self.layer5(out)
44        out = out.reshape(out.size(0), -1)
45        out = self.fc(out)
46        out = self.fc1(out)
47        out = self.fc2(out)
48    return out
49
```

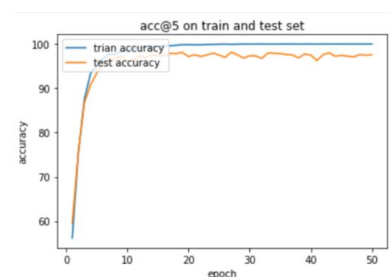
Listing 6: structure of part3



(a) loss



(b) acc@1



(c) acc@5

Figure 4

*Results :*

*maximum acc@1 on test set converges to : 76.64*  
*maximum acc@5 on test set converges to : 96*  
*maximum acc@1 on trian set converges to : 99.47*  
*maximum acc@5 on train set converges to : 100*

conclusions :

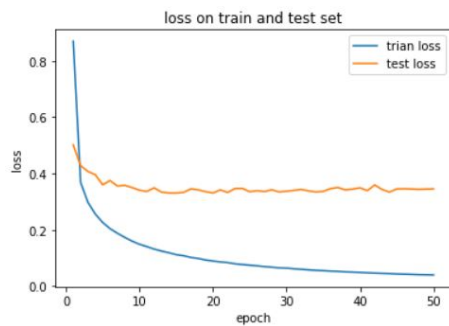
same as previous parts

## part4

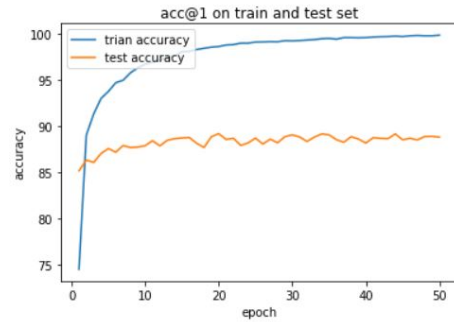
whole alexnet architecture : we squeeze the pre-trained weights ad just learn the new weight of classifier layer :

```
1 AlexNet_model = torch.hub.load('pytorch/vision:v0.6.0', 'alexnet', pretrained=True)
2 #---description
3 AlexNet_model.eval()
4 #-----
5
6
7 for param in AlexNet_model.parameters():
8     param.requires_grad = False #freeze all layers
9 for param in AlexNet_model.classifier[6].parameters():
10    param.requires_grad = True #unfreez classifer layer
```

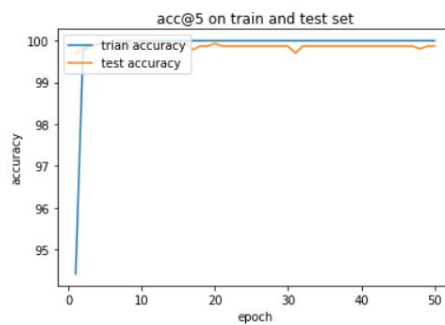
Listing 7: freeze some parameters



(a) loss



(b) acc@1



(c) acc@5

Figure 5

*Results :*

*maximum acc@1 on test set converges to : 89.2*  
*maximum acc@5 on test set converges to : 99.8*  
*maximum acc@1 on trian set converges to : 99.86*  
*maximum acc@5 on train set converges to : 100*

conclusions :

same as previous parts



## part5

transfer learning : in this part we use trained weights of alexnet model that use for imagenet dataset

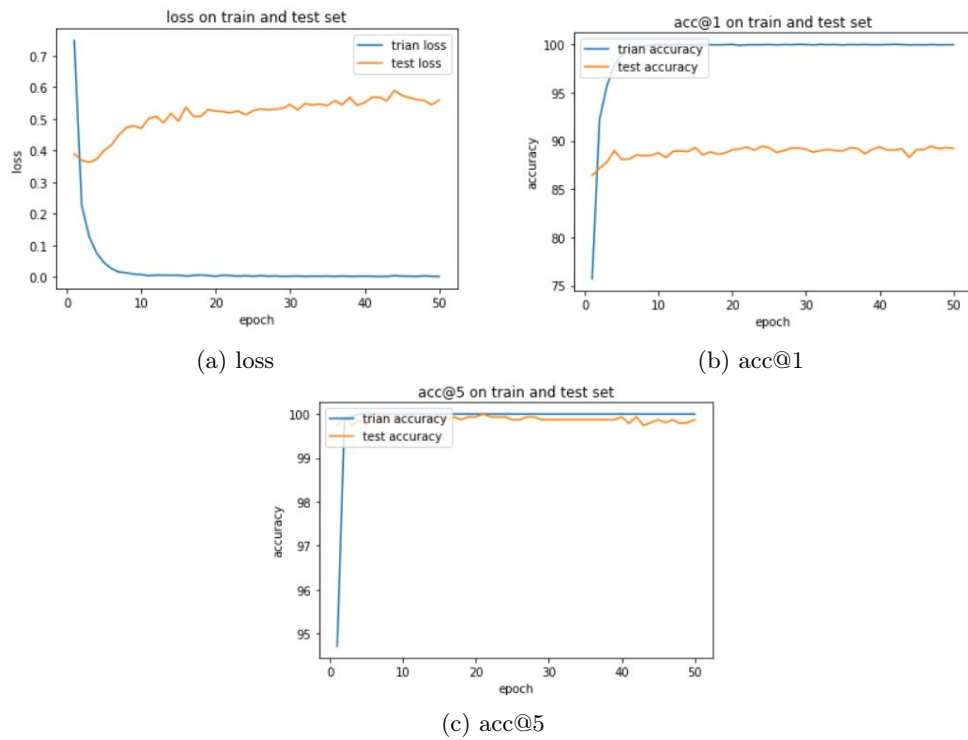


Figure 6

*Results :*

*maximum acc@1 on test set converges to : 89.43*

*maximum acc@5 on test set converges to : 99.7*

*maximum acc@1 on train set converges to : 99.93*

*maximum acc@5 on train set converges to : 100*

conclusions :

same as previous parts

## summary

results					
	optimizer	lr	acc@1	acc@5	augmentation
part1	SGD	5e-4	63.4%	93.0%	2X(horizontal flip)
part2	SGD	5e-4	73.9%	93.5%	2X(horizontal flip)
part3	SGD	5e-4	76.6%	96.0%	2X(horizontal flip)
part4	SGD	5e-4	89.2%	99.8%	2X(horizontal flip)
part5	SGD	5e-4	89.4%	99.7%	2X(horizontal flip)
* number of epoches : 50					
* number of training pictures : 5970					

Figure 7: summary

This is link of my notebook for this homework: [Link](#).