# CV HW5 Classifier

## Group 6

- 0516054 劉雨恩
- 0756049 許耕福
- 0856736 鄭鈞聯

## Introduction

This assignment is to let us know how to implenment tiny images representation, nearest neighbor classifier, bag of SIFT representation, and linear SVM classifier and use them to classify 15 catagory of images. Besides, we could also use deep learning to solve the task.

## Implementation procedure

1. Tiny images representation + nearest neighbor classifier
2. Bag of SIFT representation + nearest neighbor classifier
3. Bag of SIFT representation + linear SVM classifier
4. (Bonus) Resnet

### Tiny images representation

- In this part, we first resize image to 16 x 16, and then directly use raw pixel as the descriptor of image.
- After getting descriptor, we then can train a k-Nearest Neighbor classifier.

**Get image descriptor**

```
1    def tiny_img(img, size=(16, 16)):
2        return cv2.resize(img, size, interpolation=cv2.INTER_AREA)
3
4    imgs = []
5    for path, label in zip(train_img, train_label):
6        img = cv2.imread(os.path.join(src_path, 'train', label[0], path[0
7        img = tiny_img(img)
8        imgs.append(img)
```

## k-Nearest Neighbor classifier

- To implement knn, we first calculate the distance between testing data points and training data points.
- Next, we get the first k closest points and find the max cluster in these data points.
- To implement weighted knn, we use the inverse of distance as the weight when deciding cluster from k data points.

```
1    def knn(train_img, train_label, test_img, test_label, k):
2        correct = 0
3        for img, label in zip(test_img, test_label):
4            distances = pow(train_img - img, 2).sum(1) ** 0.5
5            index = distances.argsort()
6            neighbors = [train_label[idx] for idx in index[:k]]
7            pred_label = get_label(neighbors, distances)
8            if pred_label == label:
9                correct += 1
10        return correct * 100 / test_label.shape[0]
11
12   def get_label(neighbors, distances):
13       labels = {}
14       for label, distance in zip(neighbors, distances):
15           if label in labels:
16               labels[label] += 1. / distance
17           else:
18               labels[label] = 1. / distance
19       return max(labels.keys(), key=(lambda key: labels[key]))
```

# Bag of SIFT representation

- There will be three different ways to find SIFT representation, the first is use cv2's SIFT detection, the second is vl_sift and vl_dsift in vlfeat the third is dsift in cyvlfeat.

- There will be two return values in sift function. Key point will be a list of destination of key point and destination will be a numpy array of shape Number_of_Keypoints×128.
- Keypoint Descriptor: Now keypoint descriptor is created. A 16x16 neighbourhood around the keypoint is taken. It is divided into 16 sub-blocks of 4x4 size. For each sub-block, 8 bin orientation histogram is created. So a total of 128 bin values are available. It is represented as a vector to form keypoint descriptor. In addition to this, several measures are taken to achieve robustness against illumination changes, rotation etc.

## SIFT description in cv2

- After using this sift discription in cv2 the accuracy is around isn't good enough so we change to another approach.

```
1   sift = cv2.xfeatures2d.SIFT_create()
2   kp, des = sift.detectAndCompute(gray,None)
```

## SIFT description in vlfeat

- According to the information from internet, I found that dense sift is better in classification so i choose to use vl_dsift in vlfeat. However, this approach is using python to run matlab code which may be really slow due to the communicatin. Besides, when setting step to 10 the accurancy could be over 50%.

```
1   eng = matlab.engine.start_matlab()
2   mat_img = matlab.single(value[i][1].tolist())
3   kp, des = eng.vl_dsift(mat_img, 'step', float(10), nargout=2)
```

## Create SIFT description in dsift

- Using dsift in cyvlfeat will be the best approach. Not only the number of founded description points is really huge, but also the speed is really fast.
- Here is the code for setting step size to [5, 5] and we could find over 3 milions points.

```
1   def create_sift_discription(train_data):
2
3       description_bag = []
4       sift_data_dic = {}
5
6       for key, value in train_data.items():
7           temp_list = []
8           for i in range(len(value)):
9               kp, des = dsift(value[i][1], step=[5, 5], fast=True)
10              description_bag.extend(des)
11              temp_list.append(des)
12          sift_data_dic[key] = temp_list
13      return description_bag, sift_data_dic
```

## Run kmeans

- We use vl_kmeans to get the center and the assignment of the descriptions found in SIFT at first. Although the accuracy is great, the running speed is really slow, especially when using huge data point.
- We than choose the kmeans in cyvlfeat and the accuracy and running speed is really reasonable.

```
1   vocab = kmeans(bag_of_features, cluster_num, initialization='PLUSPLUS
```

## Create histogram

- This step is to create the histogram of every image in training data and test data. According th each images, we assign each features in the images to a cluster and add them up to from a histogram for every images.

```
1   def get_bags_of_sifts(data_dict, vocab):
2       feats, labels = [], []
3       for i, label in enumerate(data_dict.keys()):
4           for descriptors in data_dict[label]:
5               distances = distance.cdist(descriptors, vocab, 'euclidean
6               clusters = np.argmin(distances, axis=1)
7               hist, bin_edges = np.histogram(clusters, \
8                       bins=np.arange(len(vocab)+1), density=True)
9               feats.append(hist)
10              labels.append(i)
11              assert len(hist) == (len(bin_edges)-1)
12      return feats, labels
```

## Linear SVM classifier

- In this step we use libsvm to implement the linear SVM procedure. In libsvm it
  provide multiclass classification so we could call the function directly.

```
1   train_x = [{(i+1): feat[i] for i in range(len(feat))} for feat in tra
2   train_y = [(label+1) for label in train_labels]
3
4   prob = svm_problem(train_y, train_x)
5   param = svm_parameter('-t 0 -c 1e7 -b 1')
6   model = svm_train(prob, param)
7
8   test_x = [{(i+1): feat[i] for i in range(len(feat))} for feat in test
9   test_y = [(label+1) for label in test_labels]
10
11  p_label, p_acc, p_val = svm_predict(test_y, test_x, model)
```

## (Bonus) Resnet

1. Dataloader
2. Training model
3. Testing data

### 0. Hyperparameter

- Batch size: 64
- Number of Resnet layers: 50
- Loss function: nn.CrossEntropyLoss()

- Optimizer: Adam
  - weight decay: 1e-5
  - learning rate: 1e-4 (decay*0.99 per EPOCH)

## 1. Dataloader

- Resize: Since the images have different shape, we have to resize the images to the same shape as (250, 250).
- Channel: Some channels of the images are 1 (gray scale), so we have to change them into 3 channels by doing `Image.convert('RGB')`.
- Normalization: Using `torchvision.transforms.ToTensor()` can make the values of the each image be normalized in range [0, 1]. Also, transform the [H, W, C] of the image into [C, H, W]. (H: Height, W: Width, C: Channel)

```python
1   import pandas as pd
2   import torch as t
3   from torch.utils import data
4   import numpy as np
5   from PIL import Image
6   from torchvision import transforms
7
8   class dataLoader(data.Dataset):
9       def __init__(self, src_path, mode):
10          self.src_path = src_path
11          self.mode = mode
12          self.map_label = {'Bedroom': 0, 'Coast': 1, 'Forest': 2, \
13          'Highway': 3, 'Industrial': 4, 'InsideCity': 5, 'Kitchen': 6,
14          'LivingRoom': 7, 'Mountain': 8, 'Office': 9, 'OpenCountry': 1
15          'Store': 11, 'Street': 12, 'Suburb': 13, 'TallBuilding': 14}
16          self.img_name, self.label = self.getData()
17
18          self.trans = transforms.Compose([transforms.Resize((250, 250)
19                                      interpolation=Image.NEARE
20                                  transforms.ToTensor(),
21                                  transforms.Normalize( \
22                                      mean=[0.485, 0.456, 0.40
23                                      std=[0.229, 0.224, 0.225
24
25          print("> Found %d images..." % (len(self.img_name)))
26
27      def __len__(self):
28          """'return the size of dataset"""
29          return len(self.img_name)
30
31      def __getitem__(self, index):
```

```python
32          label_name = self.label[index]
33          label = self.map_label[self.label[index]]
34          path = self.src_path + self.mode \
35                  + '/' + label_name + '/' + self.img_name[index]
36
37          img_origin = Image.open(path, 'r').convert('RGB')
38          img = self.trans(img_origin)
39
40          return img, label
41
42      def getData(self):
43          img = pd.read_csv(self.src_path + self.mode + '/img.csv')
44          label = pd.read_csv(self.src_path + self.mode + '/label.csv')
45
46          return np.squeeze(img.values), np.squeeze(label.values)
47
48  # main
49  train_data = dataLoader(src_path,'train')
50  train_len = len(train_data)
51  test_data = dataLoader(src_path, 'test')
52  test_len = len(test_data)
53
54  train_loader = Data.DataLoader(
55      dataset=train_data,
56      batch_size=BATCH_SIZE,
57      shuffle=True,
58      num_workers=12,
59  )
60  test_loader = Data.DataLoader(
61      dataset=test_data,
62      batch_size=BATCH_SIZE,
63      shuffle=False,
64      num_workers=12,
65  )
66
67  for epoch_i in count(1):
68      # training
69      for batch_x, batch_y in train_loader:
70          ...
71
72      # testing
73      for batch_x, batch_y in test_loader:
74          ...
75
```

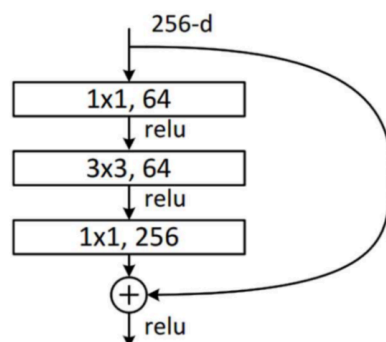## 2. Training model (Resnet 50, pretrained)

- Net arcitecture

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|---|---|---|---|---|---|---|
| conv1 | 112×112 | | | 7×7, 64, stride 2 | | |
| | | | | 3×3 max pool, stride 2 | | |
| conv2_x | 56×56 | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times23$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ |
| | 1×1 | | | average pool, 1000-d fc, softmax | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

- ResNet solves deep network problems
  - Vanishing/Exploding gradients
  - Degradation problem

- Residual Block:



Bottleneck block

$$256\text{-}d$$

```
1x1, 64
   relu
3x3, 64
   relu
1x1, 256
```

$$+$$

relu

```
 1   import torch.nn as nn
 2   import torchvision.models
 3
 4   def lr_decay(optim):
 5       for param_group in optim.param_groups:
 6           param_group['lr'] = param_group['lr']*0.99
 7
 8   class ResNet_pre(nn.Module):
 9       def __init__(self, layer, pretrained=True):
10           super(ResNet_pre, self).__init__()
11
12           if(layer==18):
13               self.classify = nn.Linear(512, 15)
14           else:
15               self.classify = nn.Linear(2048, 15)
16
17           pretrained_model = torchvision.models.__dict__ \
18                       ['resnet{}'.format(layer)](pretrained=pre
```

```python
            self.conv1 = pretrained_model._modules['conv1']
            self.bn1 = pretrained_model._modules['bn1']
            self.relu = pretrained_model._modules['relu']
            self.maxpool = pretrained_model._modules['maxpool']

            self.layer1 = pretrained_model._modules['layer1']
            self.layer2 = pretrained_model._modules['layer2']
            self.layer3 = pretrained_model._modules['layer3']
            self.layer4 = pretrained_model._modules['layer4']

            self.avgpool = nn.AdaptiveAvgPool2d(1)

        def forward(self, x):
            x = self.conv1(x)
            x = self.bn1(x)
            x = self.relu(x)
            x = self.maxpool(x)

            x = self.layer1(x)
            x = self.layer2(x)
            x = self.layer3(x)
            x = self.layer4(x)

            x = self.avgpool(x)
            x = x.view(x.size(0), -1)
            x = self.classify(x)

            return x

# main
for epoch_i in count(1):
    right = 0
    for batch_x, batch_y in train_loader:
        batch_x = batch_x.cuda()
        batch_y = batch_y.cuda()

        # clear gradient
        optimizer.zero_grad()

        # forward + backward
        output = net.forward(batch_x.float())
        y_o = torch.max(output, 1)[1]
        y_hat = sum(batch_y == y_o).item()
        right += y_hat

        loss = loss_funct(output, batch_y)
        loss.backward()

        # update parameters
```

```
68          optimizer.step()
69
70      accuracy = right/train_len
71      train_accuracy.append(accuracy * 100)
72      fp_train.write(str(accuracy * 100) + '\n')
73      print('\nEpoch {}, Train accuracy: {}'.format(epoch_i, accuracy))
74
75      ...
```

## 3. Testing data

```
1   # main
2   for epoch_i in count(1):
3       ...
4
5       # test
6       right = 0
7       net.eval()
8       with torch.no_grad():
9           for batch_x, batch_y in test_loader:
10              batch_x = batch_x.cuda()
11              batch_y = batch_y.cuda()
12              test_out = net.forward(batch_x.float())
13              test_out = torch.max(test_out, 1)[1]
14              test_hat = sum(batch_y==test_out).item()
15              right += test_hat
16      net.train()
17
18      accuracy = right / test_len
19      if(accuracy > highest_accuracy):
20          torch.save(net.state_dict(), \
21                      target_path+title+'_{}.pkl'.format(epoch_i))
22          highest_accuracy = accuracy
23          if highest_accuracy == 1:
24              break
25      test_accuracy.append(accuracy * 100)
26      fp_test.write(str(accuracy * 100) + '\n')
27      print('Test accuracy:', accuracy)
```

# Experimental results

## Tiny images representation + nearest neighbor classifier

- In this part, we found that when the k is equal to 12, the accuracy can reach to 18% in the testing dataset.

| k | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| accuracy | 13.3% | 13.3% | 13.3% | 16.7% | 16.7% |

| k | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|
| accuracy | 15.3% | 14.0% | 16.0% | 15.3% | 16.7% |

| k | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|
| accuracy | 15.3% | 18.0% | 17.3% | 17.3% | 16.7% |

## Bag of SIFT representation + nearest neighbor classifier

- There are several parameters that we could change, for example, the SIFT description, the step of vl_dsift, cluster number, and the k in knn.
- After some testing we found that using dsift in cyvlfeat, step $[5, 5]$, and set k to 1 in knn will be the most reasonable parameter. There will be some difference between the number of clusters using in kmeans.
- Ferthermore, when using knn there will be differences using norm1 or norm2, the result is the following.

norm2

| cluster num | 100 | 200 | 500 | 1000 |
|---|---|---|---|---|
| accuracy | 49.3% | 47.3% | 52.7% | 52.7% |

norm1

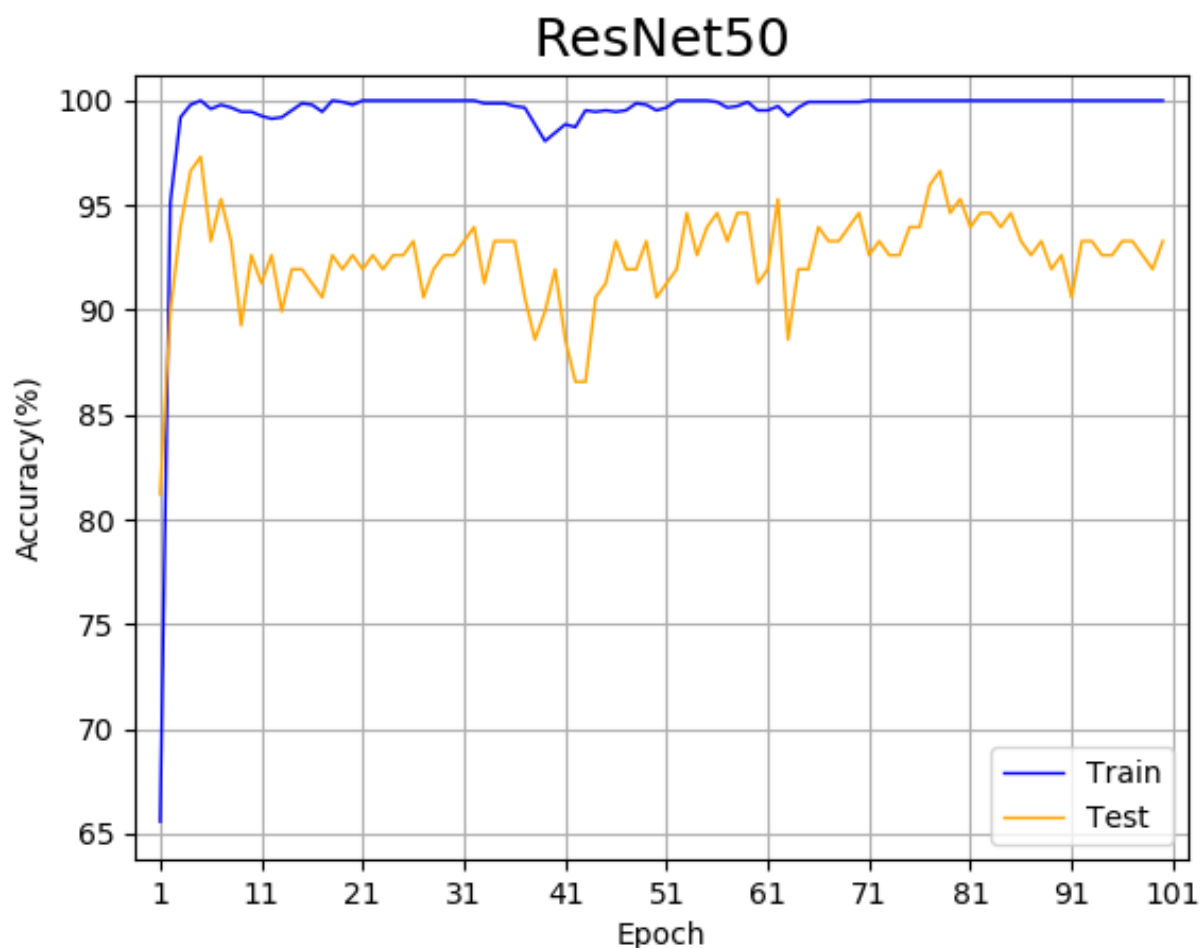| cluster num | 100 | 200 | 500 | 1000 |
|---|---|---|---|---|
| accuracy | 58% | 59.3% | 63.3% | 64.7% |

## Bag of SIFT representation + linear SVM classifier

- As same as the previous part, we could choose several SIFT description. We simply use the dsift in cyvlfeat since we did several experiments before. However, if we increase the number of cluster, it will increase the accuracy rate obviously. When

using 1000 clusters the accruacy could be 72%.

| cluster num | 100 | 200 | 500 | 1000 |
|---|---|---|---|---|
| accuracy | 63.3% | 66.7% | 68% | 72% |

## Resnet

- Training process



- Epoch: 100
  - Highest training accuracy: 100.0%
  - Highest testing accuracy: 97.315%

# Discussion

---

### nearest neighbor classifier

- In Bag of SIFT representation + nearest neighbor classifier part if we change the norm2 distance to norm1 the accurancy increased hugely. We are not sure why cause this kind of phenomenon.

## Bag of SIFT representation

- In this assingnment we are not allowed to use sklearn. However, I had met many problems using vlfeat. The first thing I met is that I couldn't setup vlfeat in python. Therefore, I tried to install vlfeat in matlab and use python to run matlab code. The result will be good enough but the running speed is really slow due to the trasmition between the two programs. Finally, I had finished setting up cyvlfeat and both the performace and speed is really reasonable.

- I had tried the sift function in cv2, vl_dsift, and dsift. The sift in cv2 could only found around 700 thousand keypoints. Using vl_dsift could found 2 milion keypoints if your step is only 5, 900 thousand when step is 10. If we use dsift with $[5,5]$ step, we could get 3 millions points. The result using vl_dsift and dsift is much better than using sift in cv2.

## linear SVM classifier

- We found that if we increase the cost in libsvm, the accuacy will increase largly. We use cost 1 at first with 200 clusters, the accuracy is only 38%. If we increase it to 4, it will turn to 44%. Last, we change it to $10^7$ than the accracy will be 66.7%.

- The reason why adding the cost will increase the performace is in SVM we would like to minimum the function $\frac{1}{2}\|w\|^2 + C\sum_i \xi_i$ subject to $y_i\left(w^T x_i - b\right) - 1 + \xi_i \geq 0 \quad \forall i, \xi_i \geq 0 \quad \forall i$. There is a cost $C$ in this function which represent the penalty of error terms $\xi_i$ If we do not give the enough penaty the $\frac{1}{2}\|w\|^2$ part will be useless.

## Resnet

- If the `shuffle` of the `Data.DataLoader` isn't `True`, the testing accuracy would be very bad as 3%. It's because the model overfits each class every time during each batch, and finally can't learn general case.

- As the batch size is increased, the accuracy also improves a lot, but the memory usage will change relatively. As a result, considering the machine performance, the improvement of the batch size will be limited.

- When the weight decay mechanism is added, the testing accuracy improved a little bit, while the training accuracy dropped a little bit. It's because the weight decay mechanism helps the model not to overfit the training data too much.

# Conclusion

After finishing the assignment, we know about several machine learning tecniques to solve classification problem in computer vision. Besides, we also learn about how to use deep learning tecnique such as Resnet to solve the classification problems.

# Work assignment plan

- We wrote our own code separately.
    - Bag of SIFT representation, nearest neighbor classifier & linear SVM: 0856736 鄭鈞聯
    - Tiny images representation & nearest neighbor classifier: 0756049 許耕福
    - Resnet: 0516054 劉雨恩
- Then, we use Messenger to discuss the problem we had encountered and finished the report together.