

# CS 671 Homework 4

Yutong Shao  
NetID ys392

December 9, 2022

**I consent to the following agreements.**

*This assignment represents my own work. I did not work on this assignment with others.  
All coding was done by myself.*

*I understand that if I struggle with this assignment that I will reevaluate whether this is the  
correct class for me to take. I understand that the homework only gets harder.*

## 1 VC dimension of Binary Decision Trees with Fixed Split Points

*Proof.* Proof by induction.

If  $\mathcal{F} := \{\text{the set of all binary decision trees whose split points are exactly equal to } \mathcal{X} \text{ with exactly 2 leaves}\}$ , then this class of trees can at most classify 2 points, because we can write four labeling assignments for this decision tree:  $(1, 1), (1, -1), (-1, 1), (-1, -1)$ .

But three points can not be shattered by this class of trees. We can prove this by contradiction. Assume we can use  $F \in \mathcal{F}$  to shatter three points, then there must exist two points that are assigned the same label. However, these two data points are not necessarily in the same class, which contradicts with the assumption.

By induction, for  $\mathcal{F} := \{\text{the set of all binary decision trees whose split points are exactly equal to } \mathcal{X} \text{ with exactly } \ell \text{ leaves}\}$ , then this class of trees can at most classify  $\ell$  points. And Similarly, we can prove by contradiction that  $\ell + 1$  points cannot be perfectly classified.

Therefore,  $VC(\mathcal{F}) = \ell$ .

□

## 2 Topic Modeling with EM

### 2.1 Derive Log-Likelihood

$\therefore$  Word  $n$  are drawn i.i.d. from topic  $k$

$\therefore$  The total probability that word  $n$  appear in document  $i$  is

$$\text{likelihood}(\theta) = \prod_{q=1}^Q \sum_{k=1}^K \Pr(w_n|z_k) \Pr(z_k|d_i)$$

where  $Q$  is the total number of words,  $Q = \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i)$ .

$\therefore$  The log likelihood is

$$\begin{aligned} \log \text{likelihood}(\theta) &= \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \log \sum_{k=1}^K \Pr(w_n|z_k) \Pr(z_k|d_i) \\ &= \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \log \sum_{k=1}^K \beta_{kn} \alpha_{ik} \end{aligned}$$

### 2.2 E Step

$$\begin{aligned} p(z_k|d_i, w_n, \alpha^{old}, \beta^{old}) &= \frac{\Pr(w_n|z_k, d_i, \alpha^{old}, \beta^{old}) \Pr(z_k|d_i, \alpha^{old}, \beta^{old})}{\Pr(w_n|\alpha^{old}, \beta^{old})} \\ &= \frac{\alpha^{old} \beta^{old}}{\sum_{k=1}^K \alpha^{old} \beta^{old}} \end{aligned}$$

### 2.3 Find ELBO for M-Step

$$\begin{aligned} \log \text{likelihood}(\theta) &= \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \log \sum_{k=1}^K \Pr(w_n|z_k) \Pr(z_k|d_i) \\ &= \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \log \sum_{k=1}^K \gamma_{ink} \frac{\beta_{kn} \alpha_{ik}}{\gamma_{ink}} \\ &\geq \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log \frac{\beta_{kn} \alpha_{ik}}{\gamma_{ink}} \quad \leftarrow \text{Jensen's inequality} \\ &:= A(\alpha, \beta) \end{aligned}$$

## 2.4 M Step

Define

$$A'(\alpha, \beta) := \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log \beta_{kn} \alpha_{ik}$$

since when maximizing auxiliary function, we can only focus on numerator.

Therefore, the constrained maximization problem is

$$\begin{aligned} & \max_{\alpha, \beta} A'(\alpha, \beta) \\ \text{s.t. } & \sum_{k=1}^K \alpha_{ik} = 1 \\ & \sum_{n=1}^N \beta_{kn} = 1 \end{aligned}$$

Write Lagrangian as

$$\mathcal{L}(\alpha, \beta) = A'(\alpha, \beta) + \lambda_1 \left( 1 - \sum_{k=1}^K \alpha_{ik} \right) + \lambda_2 \left( 1 - \sum_{n=1}^N \beta_{kn} \right)$$

Compute derivatives and set to 0.

$$\frac{\partial \mathcal{L}}{\partial \alpha_{ik}} = 0 \tag{1}$$

$$\frac{\partial \mathcal{L}}{\partial \beta_{kn}} = 0 \tag{2}$$

For derivative (1),

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \alpha_{ik}} &= \frac{\partial A'(\alpha, \beta)}{\partial \alpha_{ik}} - \lambda_1 \\ &= \frac{\partial}{\partial \alpha_{ik}} \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log \beta_{kn} \alpha_{ik} - \lambda_1 \\ &= \frac{\partial}{\partial \alpha_{ik}} \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log (\beta_{kn} + \alpha_{ik}) - \lambda_1 \\ &= \underbrace{\frac{\partial}{\partial \alpha_{ik}}}_{\text{derivative of certain i and k}} \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log \alpha_{ik} - \lambda_1 \\ &= \sum_{n=1}^N q(w_n; d_i) \gamma_{ink} \frac{1}{\alpha_{ik}} - \lambda_1 = 0 \\ \implies \alpha^{new} &= \frac{\sum_{n=1}^N q(w_n; d_i) \gamma_{ink}}{\lambda_1} \end{aligned}$$

Note that  $\sum_{k=1}^K \alpha_{ik} = 1$ . Thus,

$$\begin{aligned}
 \sum_{k=1}^K \alpha_{ik} &= \sum_{k=1}^K \frac{\sum_{n=1}^N q(w_n; d_i) \gamma_{ink}}{\lambda_1} = 1 \\
 \implies \lambda_1 &= \sum_{k=1}^K \sum_{n=1}^N q(w_n; d_i) \gamma_{ink} \\
 &= \sum_{n=1}^N q(w_n; d_i) \underbrace{\sum_{k=1}^K \gamma_{ink}}_{=1} \\
 &= \sum_{n=1}^N q(w_n; d_i)
 \end{aligned}$$

Therefore,

$$\alpha^{new} = \frac{\sum_{n=1}^N q(w_n; d_i) \gamma_{ink}}{\sum_{n=1}^N q(w_n; d_i)}$$

We can also derive  $\beta^{new}$  follow the same logic. According to (2),

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \beta_{kn}} &= \underbrace{\frac{\partial}{\partial \beta_{kn}}}_{\text{derivative of certain k and n}} \sum_{i=1}^M \sum_{n=1}^N \sum_{k=1}^K q(w_n; d_i) \gamma_{ink} \log \beta_{kn} - \lambda_2 \\
 &= \sum_{i=1}^M q(w_n; d_i) \gamma_{ink} \frac{1}{\beta_{kn}} - \lambda_2 = 0 \\
 \implies \beta^{new} &= \frac{\sum_{i=1}^M q(w_n; d_i) \gamma_{ink}}{\lambda_2}
 \end{aligned}$$

Note that  $\sum_{n=1}^N \beta_{kn} = 1$ . Thus,

$$\begin{aligned}
 \sum_{n=1}^N \beta_{kn} &= \sum_{n=1}^N \frac{\sum_{i=1}^M q(w_n; d_i) \gamma_{ink}}{\lambda_2} = 1 \\
 \implies \lambda_2 &= \sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \gamma_{ink}
 \end{aligned}$$

Therefore,

$$\beta^{new} = \frac{\sum_{i=1}^M q(w_n; d_i) \gamma_{ink}}{\sum_{i=1}^M \sum_{n=1}^N q(w_n; d_i) \gamma_{ink}}$$

### 3 Gradient computations in Neural Networks

#### 3.1 Gradient Evaluation

Define the NN's structure as follows.

Input layer ( $l = 0$ ) - hidden layer 1 ( $l = 1$ ) - hidden layer 2 ( $l = 2$ ) - Output Layer ( $l = 3$ )

Before diving into derivations, I want to change some notations to make the vector and matrix representation clearer.

Let  $\mathbf{X}$  be the input vector,  $\mathbf{X} = (x_1, x_2, \dots, x_d)_{(d \times 1)}^T$ ;

Let  $\mathbf{W}^{(l)}$  be the weight matrix of layer  $l$ . Specifically,  $\mathbf{W}^{(1)}$  should be a  $d \times H$  matrix,  $\mathbf{W}^{(2)}$  is  $H \times H$ ,  $\mathbf{W}^{(3)}$  should be  $H \times 1$ . For example,

$$\mathbf{W}^{(2)} = \begin{pmatrix} w_{11}^{(2)} & w_{11}^{(2)} & \dots & w_{1H}^{(2)} \\ w_{21}^{(2)} & w_{21}^{(2)} & \dots & w_{2H}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{H1}^{(2)} & w_{H1}^{(2)} & \dots & w_{HH}^{(2)} \end{pmatrix}$$

Let  $\mathbf{b}^{(l)}$  be the bias vector of layer  $l$ , which should be  $H \times 1$  for layer 1 and 2, and a scalar for  $l = 3$ ;

$$\begin{aligned} \mathbf{b}^{(l)} &= \left( b_1^{(l)}, b_2^{(l)}, \dots, b_H^{(l)} \right)_{(H \times 1)}^T, \quad l = 1, 2 \\ \mathbf{b}^{(l)} &= b^{(3)}, \quad l = 3 \end{aligned}$$

Let  $\mathbf{Z}^{(l)}$  denote the pre-activation vector at layer  $l$ ,

$$\begin{aligned} \mathbf{Z}^{(l)} &= \mathbf{W}^{(l)T} \mathbf{X} + \mathbf{b}^{(l)}, \quad l = 1 \\ \mathbf{Z}^{(l)} &= \mathbf{W}^{(l)T} \mathbf{h}^{(l-1)} + \mathbf{b}^{(l)}, \quad l > 1 \end{aligned}$$

And,  $\mathbf{h}^{(l)}$  denote the activation function vector of layer  $l$ .

$$\begin{aligned} \mathbf{h}^{(l)} &= \sigma \left( \mathbf{Z}^{(l)} \right) = \left( h_1^{(l)}, h_2^{(l)}, \dots, h_H^{(l)} \right)_{(H \times 1)}^T \\ &= \left( \sigma \left( z_1^{(l)} \right), \sigma \left( z_2^{(l)} \right), \dots, \sigma \left( z_H^{(l)} \right) \right)_{(H \times 1)}^T \end{aligned}$$

Finally, the derivative of sigmoid function is

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

1. From Output Layer ( $l = 3$ ) to hidden layer 2 ( $l = 2$ )

Goal: compute  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(3)}}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(3)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(3)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{Z}^{(3)}} \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(3)}}$$

where

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(3)}} &= -\frac{y}{\mathbf{h}^{(3)}} + \frac{1-y}{1-\mathbf{h}^{(3)}} \\ \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{Z}^{(3)}} &= \mathbf{h}^{(3)} (1 - \mathbf{h}^{(3)}) \\ \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(3)}} &= \left( h_1^{(2)}, h_2^{(2)}, \dots, h_H^{(2)} \right)_{(H \times 1)}^T = \mathbf{h}^{(2)} \end{aligned}$$

Therefore,

$$\begin{aligned} \delta_3 &= \mathbf{h}^{(3)} - y \\ \Rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(3)}} &= \left( \mathbf{h}^{(3)} - y \right) \mathbf{h}^{(2)} = \delta_3 \mathbf{h}^{(2)} \quad \leftarrow \text{a } 1 \times H \text{ vector} \end{aligned}$$

2. From hidden layer 2 ( $l = 2$ ) to hidden layer 1 ( $l = 1$ ).

Goal: compute  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(2)}}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(2)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{Z}^{(3)}} \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{Z}^{(2)}} \frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{W}^{(2)}}$$

where

$$\delta_2 = \frac{\partial \mathcal{L}}{\partial \mathbf{Z}^{(2)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{Z}^{(3)}} \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{Z}^{(2)}} = \delta_3 \mathbf{W}_{H \times 1}^{(3)} \underbrace{\mathbf{h}^{(2)} (1 - \mathbf{h}^{(2)})}_{1 \times H}$$

Thus,  $\frac{\partial \mathcal{L}}{\partial \mathbf{Z}^{(3)}} \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{Z}^{(2)}}$  is a  $H \times H$  matrix.

Now we left one last term  $\frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{W}^{(2)}}$ , which is a little bit complicated because we need to compute the derivative of a vector over a matrix. This should yield a high-dimension tensor, but note that we can write it as a matrix (See the math below 3.2 for the derivation). Also, we only care about the final product, which should also be a  $H \times H$  matrix.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(2)}} = \delta_3 \mathbf{W}_{H \times 1}^{(3)} \underbrace{\mathbf{h}^{(2)} (1 - \mathbf{h}^{(2)})}_{1 \times H} \mathbf{h}^{(1)T}$$

3. From hidden layer 1 ( $l = 1$ ) to the input layer ( $l = 0$ )

Goal: compute  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}}$

Similarly, we can get the derivative follow the logic above. The final answer should be a  $H \times H$  matrix.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{Z}^{(1)}} \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} \\ &= \delta_2 \underbrace{\mathbf{W}^{(2)}}_{H \times H} \underbrace{\mathbf{h}^{(1)} (1 - \mathbf{h}^{(1)})}_{H \times 1} \underbrace{X^T}_{1 \times H} \end{aligned}$$

### 3.2 General Rule

Since  $1 \leq l \leq L - 1$ , we can easily get the general form from the above reasoning.

$$\delta_l = \delta_{l+1} \mathbf{W}_{l+1} \mathbf{h}_l (1 - \mathbf{h}_l)$$

and  $\delta_l^i$  is the  $i$ -th element of vector  $\delta_l$ .

### Derivative of vector over matrix

Without the loss of generality, consider  $\mathbf{y} = \mathbf{x}W$ , where  $\mathbf{y} = [y_1, y_2, \dots, y_H]$  is a  $1 \times H$  vector, and  $\mathbf{x} = [x_1, x_2, \dots, x_d]$  is a  $1 \times d$  vector,  $W$  is a  $d \times H$  matrix.

We want to find  $\frac{d\mathbf{y}}{dW}$ . Normally, the derivative is a tensor. But we can actually write it out as a matrix.

Rewrite  $\mathbf{y}$  as elements so that the  $j$ -th element can be represented as

$$y_j = x_1 W_{1,j} + x_2 W_{2,j} + \dots + x_d W_{d,j}$$

Therefore,

$$\frac{\partial y_j}{\partial W_{i,j}} = x_i$$

Note that the derivative of  $y_j$  over other elements is 0. Thus,

$$\frac{d\mathbf{y}}{dW} = \begin{pmatrix} \frac{\partial y_1}{\partial W_{1,1}} & \frac{\partial y_1}{\partial W_{2,1}} & \frac{\partial y_1}{\partial W_{3,1}} & \dots & \frac{\partial y_1}{\partial W_{d,1}} \\ \frac{\partial y_2}{\partial W_{1,2}} & \frac{\partial y_2}{\partial W_{2,2}} & \frac{\partial y_2}{\partial W_{3,2}} & \dots & \frac{\partial y_2}{\partial W_{d,2}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_H}{\partial W_{1,H}} & \frac{\partial y_H}{\partial W_{2,H}} & \frac{\partial y_H}{\partial W_{3,H}} & \dots & \frac{\partial y_H}{\partial W_{d,H}} \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & x_3 & \dots & x_d \\ x_1 & x_2 & x_3 & \dots & x_d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1 & x_2 & x_3 & \dots & x_d \end{pmatrix}_{H \times d}$$

## 4 Clustering

See coding appendix.



## 5 Convolutional Neural Network on CIFAR-10

See coding appendix.

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import random
5 import time
6
7 import warnings
8 warnings.filterwarnings('ignore')

```

# Q4 Clustering

## 4.1 load dataset and draw scatter plot

```

1 df5 = pd.read_csv('mall_customers.csv')
2 df5.head()

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

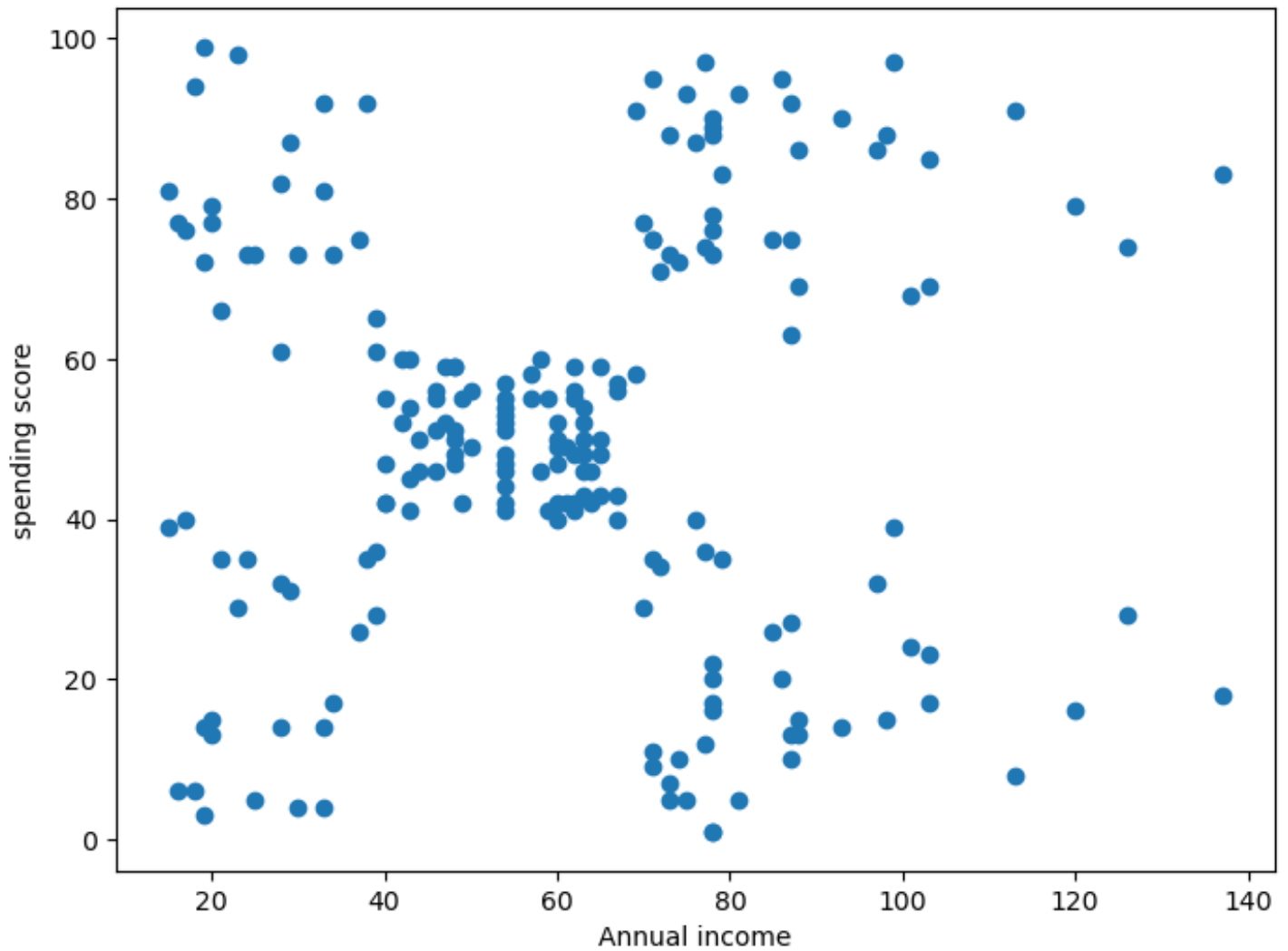
```

	Annual_Income	Spending_Score
0	15	39
1	15	81
2	16	6
3	16	77
4	17	40

```

1 fig1, ax1 = plt.subplots(figsize=(8,6))
2 plt.scatter(df5.iloc[:,0], df5.iloc[:,1]);
3 plt.xlabel('Annual income')
4 plt.ylabel('spending score');

```



## 4.2 K-means algorithm without sklearn

```

1 def EuclideanDistance(data, centers, k):
2     """
3     This function calculates the distance between each
4     data point and each center, stores in a matrix with
5     n rows (n data points) and k columns (k centers)
6     """
7     dis_mat = [[0]*k for i in range(data.shape[0])]
8     for i in range(data.shape[0]):

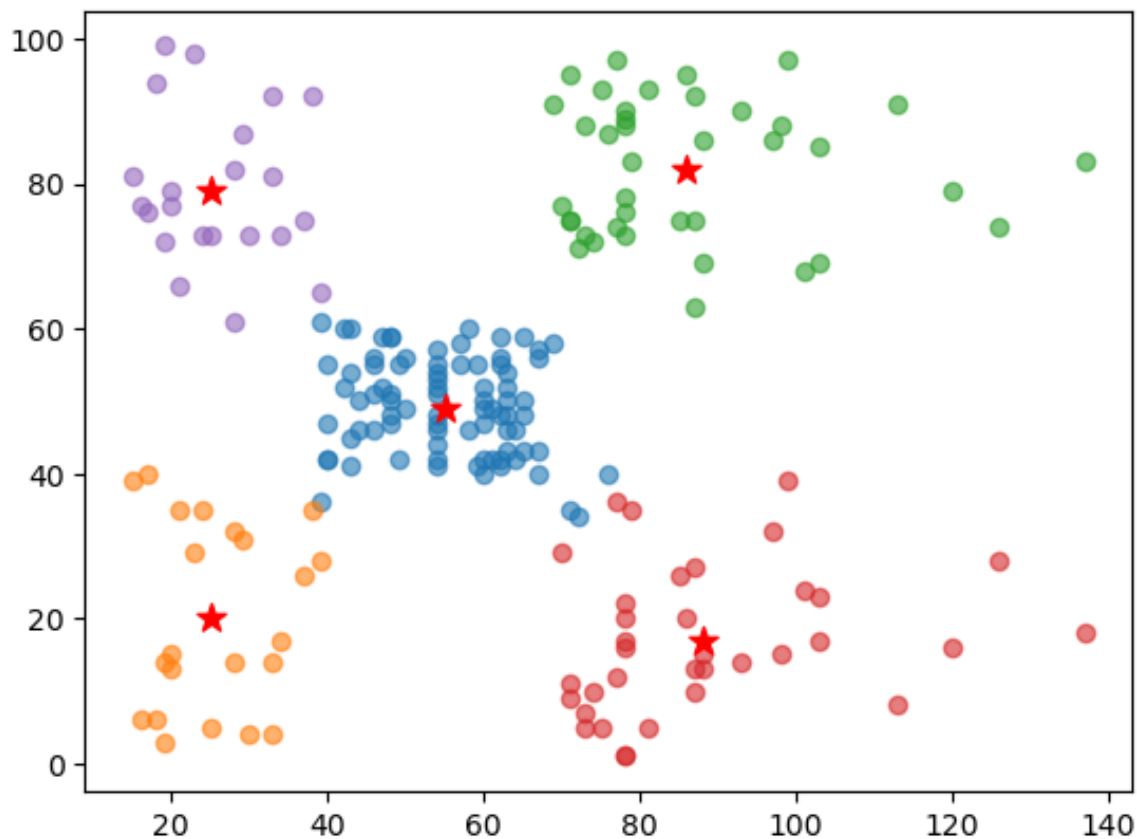
```

```

9         for j in range(k):
10             dis_mat[i][j] = np.sqrt(np.sum((data.iloc[i,:]-centers[j])**2))
11
12     return dis_mat
13
14
15 data = df5
16 epochs = 10
17 k = 5
18 rng = np.random.default_rng()
19 Cs = rng.choice(data, k)
20
21 for e in range(epochs):
22     distance = EuclideanDistance(data=df5, centers=Cs, k=5)
23     nearest_c = np.argmin(distance, axis=1)
24     for new_c in range(k):
25         Cs[new_c] = data.loc[nearest_c == new_c].mean()
26
27 for i in range(k):
28     cluster = data.loc[nearest_c == i]
29     X = cluster.iloc[:,0]
30     y = cluster.iloc[:,1]
31     plt.scatter(X, y, c=colors[i], alpha=0.6)
32
33 plt.scatter(Cs[:,0], Cs[:,1], s=100, marker='*', c='r');
34
35 print(Cs, '\n', nearest_c)

```

[illegible]



Each cluster represents:

- the upper right cluster: people who have high income and also spend much
- the lower right cluster: people who have relatively low income but spend much
- the center cluster: people with medium income and medium spending score
- the upper left cluster: people with high income but low spending score
- the lower left cluster: people with low income and low spending score

```

1 # print(Cs, '\n', nearest_c)
2 def PlotCluster(data, centers, nearest_c, k, colors):
3     fig52, ax52 = plt.subplots(figsize=(8,6))
4     plt.scatter(centers[:,0], centers[:,1], s=100, marker='*',c='r')
5     #     color = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728', '#9467bd']
6
7     for i in range(k):
8         cluster = data.loc[nearest_c == i]
9         X = cluster.iloc[:,0]
10        y = cluster.iloc[:,1]
11        plt.scatter(X, y, c=colors[i], alpha=0.6)
12
13 colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728', '#9467bd']
14 # PlotCluster(df5, Cs, nearest_c, 5, colors)

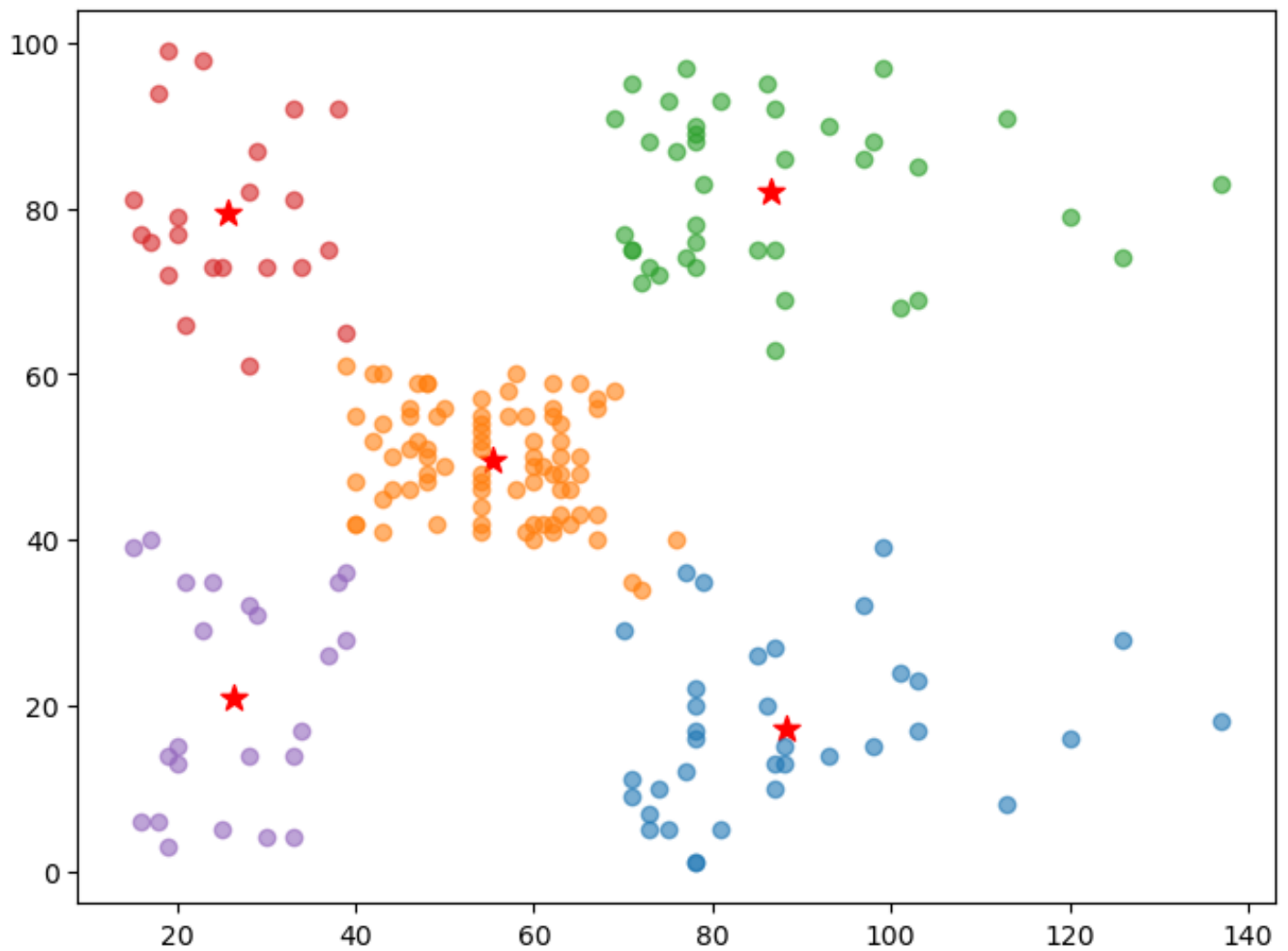
```

### 4.3 K-means with sklearn

```
1 | from sklearn.cluster import KMeans
```

```
1 kmeans = KMeans(n_clusters=5, random_state=0).fit(df5)
2 centers_sk = kmeans.cluster_centers_
3 nearest_c_sk = kmeans.labels_
4 print(centers_sk, '\n', nearest_c_sk)
5 PlotCluster(df5, centers_sk, nearest_c_sk, 5, colors);
6
```

[illegible]

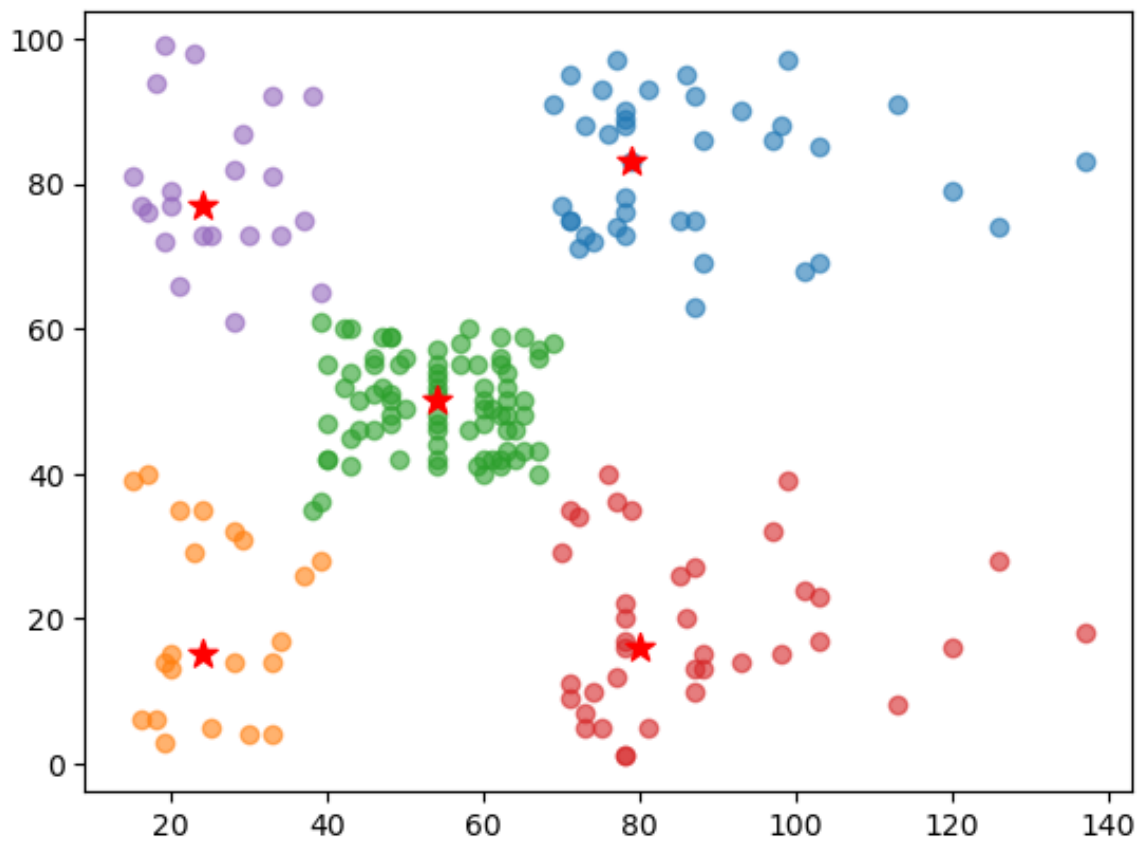


Yes, the results are the same for over 80% of the time. But Kmeans are more stable because the cluster assignment are always the same, but my algorithm has different labels for some iteration. This is probably because the initial centers of my algorithm are randomly chosen at the beginning, which may sometimes influence the final cluster result in the end.

## 4.4 K-medians







The results are not the same when using k-means and k-medians.

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torchvision
5 import torchvision.transforms as transforms
6 import torch.optim as optim
7 import matplotlib.pyplot as plt
8 import numpy as np
9 from torch.utils.data.dataset import random_split
10 from torchvision import datasets
11 from sklearn.metrics import confusion_matrix
12 import PIL
```

## Q5

```
1 ##Do Not Touch This Cell
2
3 class Net(nn.Module):
4     def __init__(self):
5         super(Net, self).__init__()
6         self.conv1 = nn.Conv2d(3, 8, 5)
7         self.conv2 = nn.Conv2d(8, 16, 3)
8         self.bn1 = nn.BatchNorm2d(8)
9         self.bn2 = nn.BatchNorm2d(16)
10        self.fc1 = nn.Linear(16*6*6, 120)
11        self.fc2 = nn.Linear(120, 84)
12        self.fc3 = nn.Linear(84, 10)
13
14        def forward(self, x):
15            out = F.relu(self.bn1(self.conv1(x)))
16            out = F.max_pool2d(out, 2)
17            out = F.relu(self.bn2(self.conv2(out)))
18            out = F.max_pool2d(out, 2)
19            out = out.view(out.size(0), -1)
20            out = F.relu(self.fc1(out))
21            out = F.relu(self.fc2(out))
22            out = self.fc3(out)
23            return out
24
25
```

```
1  ##Do Not Touch This Cell
2
3  device = 'cuda' if torch.cuda.is_available() else 'cpu'
4  net = Net().to(device)
5  optimizer = optim.SGD(net.parameters(), lr=0.01, momentum=0.5)
6  if device == 'cuda':
7      print("Train on GPU...")
8  else:
9      print("Train on CPU...")
```

```
1  Train on GPU...
```

```
1  ##Do Not Touch This Cell
2  max_epochs = 50
3
4  random_seed = 671
5  torch.manual_seed(random_seed)
```

```
1  <torch._C.Generator at 0x7f6528454170>
```

```
1  train_transform = transforms.Compose(
2      [transforms.ToTensor(),
3       transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))])
4
5  test_transform = transforms.Compose(
6      [transforms.ToTensor(),
7       transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))])
8
9  dataset = torchvision.datasets.CIFAR10(root='./data', train=True, download=True,
10                                         transform=train_transform)
11  ##TODO: Split the set into 80% train, 20% validation (there are 50K total images)
12  train_num = int(0.8 * 50000)
13  val_num = int(0.2 * 50000)
14  train_set, val_set = random_split(dataset, [train_num, val_num])
15
16  train_loader = torch.utils.data.DataLoader(train_set, batch_size=128, shuffle=True)
17  val_loader = torch.utils.data.DataLoader(val_set, batch_size=128, shuffle=False)
18
19  test_set = torchvision.datasets.CIFAR10(root='./data', train=False,
20                                           download=True, transform=test_transform)
```

```
21 test_loader = torch.utils.data.DataLoader(test_set, batch_size=1, shuffle=False)
22
23 classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog',
24            'frog', 'horse', 'ship', 'truck']
```

```
1 Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to ./data/cifar-10-python.tar.gz
```

```
1 0%|          | 0/170498071 [00:00<?, ?it/s]
```

```
1 Extracting ./data/cifar-10-python.tar.gz to ./data
2 Files already downloaded and verified
```

```
1 len(train_set)
```

```
1 40000
```

```
1 loss_list, acc_list = [], []
2 loss_list_val, acc_list_val = [], []
3 criterion = nn.CrossEntropyLoss()
4
5 for epoch in range(max_epochs):
6     #TODO: set the net to train mode:
7     model_t = net.train()
8
9     epoch_loss = 0.0
10    correct = 0
11    for batch_idx, (data, labels) in enumerate(train_loader):
12        data, labels = data.to(device), labels.to(device)
13
14        optimizer.zero_grad()
15        ##TODO: pass the data into the network and store the output
16        output = model_t(data)
17
18        ##TODO: Calculate the cross entropy loss between the output and target
19        loss = criterion(output, labels)
```

```

20
21     ##TODO: Perform backpropagation
22     optimizer.zero_grad()
23     loss.backward()
24     optimizer.step()
25
26     ##TODO: Get the prediction from the output
27     _, predicted = torch.max(output.data, 1)
28
29     ##TODO: Calculate the correct number and add the number to correct
30     correct += (predicted == labels).sum()
31
32     ##TODO: Add the loss to epoch_loss.
33     epoch_loss += loss
34
35     ##TODO: calculate the average loss
36     avg_loss = epoch_loss / len(train_set)
37
38     ##TODO: calculate the average accuracy
39     avg_acc = correct / len(train_set)
40
41     ##TODO: append average epoch loss to loss list
42     loss_list.append(avg_loss)
43
44     ##TODO: append average accuracy to accuracy list
45     acc_list.append(avg_acc)
46
47
48
49     # validation
50     ##TODO: set the model to eval mode
51     model_val = net.eval()
52
53     with torch.no_grad():
54         loss_val = 0.0
55         correct_val = 0
56         for batch_idx, (data, labels) in enumerate(val_loader):
57             data, labels = data.to(device), labels.to(device)
58             ##TODO: pass the data into the network and store the output
59             output_val = model_val(data)
60
61             ##TODO: Calculate the cross entropy loss between the output and target
62             loss_val_ = criterion(output_val, labels)
63
64             ##TODO: Get the prediction from the output
65             _, predicted_val = torch.max(output_val.data, 1)
66
67             ##TODO: Calculate the correct number and add the number to correct_val
68             correct_val += (predicted_val == labels).sum()

```

```

69
70     ##TODO: Add the loss to loss_val
71     loss_val += loss_val_
72
73     ##TODO: calculate the average loss of validation
74     avg_loss_val = loss_val / len(test_set)
75
76     ##TODO: calculate the average accuracy of validation
77     avg_acc_val = correct_val / len(test_set)
78
79     ##TODO: append average epoch loss to loss list of validation
80     loss_list_val.append(avg_loss_val)
81
82     ##TODO: append average accuracy to accuracy list of validation
83     acc_list_val.append(avg_acc_val)
84
85     print('[epoch %d] loss: %.5f accuracy: %.4f val loss: %.5f val accuracy: %.4f' %
    (epoch + 1, avg_loss, avg_acc, avg_loss_val, avg_acc_val))

```

```

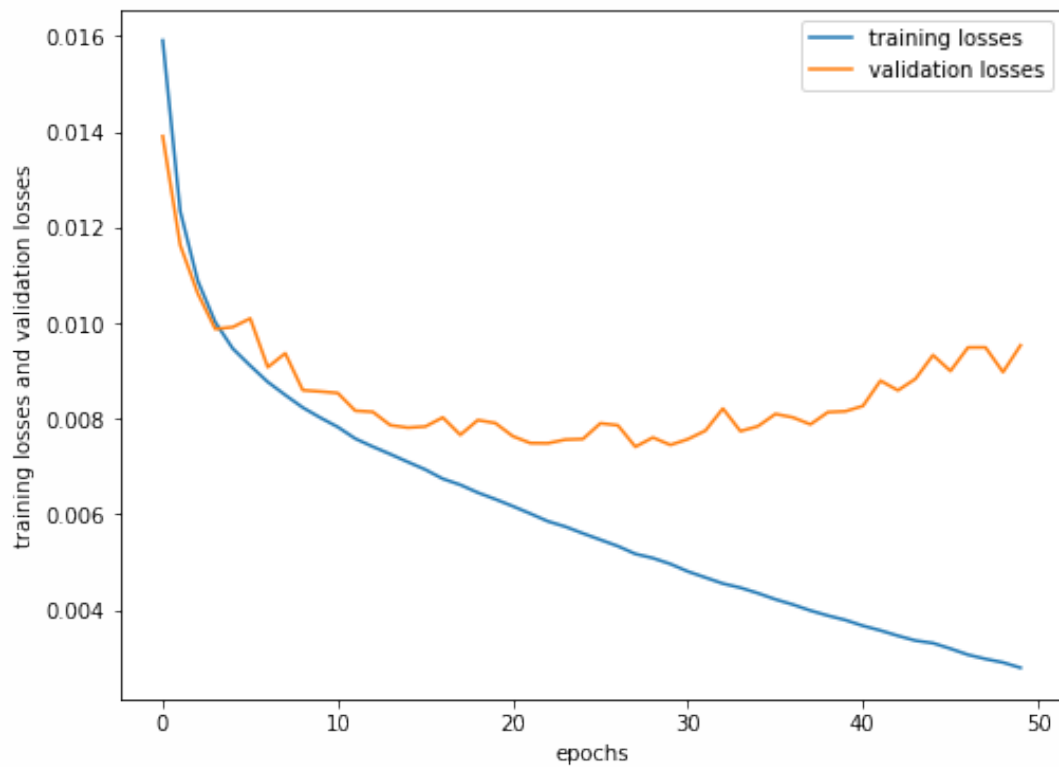
1 [epoch 1] loss: 0.01591 accuracy: 0.2610 val loss: 0.01390 val accuracy: 0.3579
2 [epoch 2] loss: 0.01234 accuracy: 0.4176 val loss: 0.01163 val accuracy: 0.4683
3 [epoch 3] loss: 0.01088 accuracy: 0.4937 val loss: 0.01062 val accuracy: 0.5151
4 [epoch 4] loss: 0.01002 accuracy: 0.5396 val loss: 0.00987 val accuracy: 0.5547
5 [epoch 5] loss: 0.00947 accuracy: 0.5695 val loss: 0.00991 val accuracy: 0.5523
6 [epoch 6] loss: 0.00911 accuracy: 0.5879 val loss: 0.01009 val accuracy: 0.5614
7 [epoch 7] loss: 0.00877 accuracy: 0.6035 val loss: 0.00908 val accuracy: 0.5976
8 [epoch 8] loss: 0.00850 accuracy: 0.6174 val loss: 0.00937 val accuracy: 0.5924
9 [epoch 9] loss: 0.00823 accuracy: 0.6284 val loss: 0.00859 val accuracy: 0.6161
10 [epoch 10] loss: 0.00802 accuracy: 0.6374 val loss: 0.00857 val accuracy: 0.6217
11 [epoch 11] loss: 0.00783 accuracy: 0.6459 val loss: 0.00854 val accuracy: 0.6227
12 [epoch 12] loss: 0.00758 accuracy: 0.6597 val loss: 0.00817 val accuracy: 0.6356
13 [epoch 13] loss: 0.00742 accuracy: 0.6677 val loss: 0.00814 val accuracy: 0.6411
14 [epoch 14] loss: 0.00726 accuracy: 0.6743 val loss: 0.00786 val accuracy: 0.6557
15 [epoch 15] loss: 0.00710 accuracy: 0.6805 val loss: 0.00781 val accuracy: 0.6534
16 [epoch 16] loss: 0.00694 accuracy: 0.6892 val loss: 0.00784 val accuracy: 0.6517
17 [epoch 17] loss: 0.00674 accuracy: 0.6996 val loss: 0.00803 val accuracy: 0.6459
18 [epoch 18] loss: 0.00662 accuracy: 0.7035 val loss: 0.00766 val accuracy: 0.6635
19 [epoch 19] loss: 0.00646 accuracy: 0.7102 val loss: 0.00797 val accuracy: 0.6535
20 [epoch 20] loss: 0.00632 accuracy: 0.7159 val loss: 0.00791 val accuracy: 0.6545
21 [epoch 21] loss: 0.00617 accuracy: 0.7240 val loss: 0.00763 val accuracy: 0.6642
22 [epoch 22] loss: 0.00602 accuracy: 0.7325 val loss: 0.00749 val accuracy: 0.6730
23 [epoch 23] loss: 0.00586 accuracy: 0.7392 val loss: 0.00749 val accuracy: 0.6749
24 [epoch 24] loss: 0.00574 accuracy: 0.7425 val loss: 0.00756 val accuracy: 0.6701
25 [epoch 25] loss: 0.00560 accuracy: 0.7512 val loss: 0.00757 val accuracy: 0.6788
26 [epoch 26] loss: 0.00547 accuracy: 0.7553 val loss: 0.00790 val accuracy: 0.6628
27 [epoch 27] loss: 0.00534 accuracy: 0.7616 val loss: 0.00786 val accuracy: 0.6643
28 [epoch 28] loss: 0.00517 accuracy: 0.7694 val loss: 0.00741 val accuracy: 0.6842
29 [epoch 29] loss: 0.00509 accuracy: 0.7736 val loss: 0.00761 val accuracy: 0.6747

```

```
30 [epoch 30] loss: 0.00496 accuracy: 0.7781 val loss: 0.00745 val accuracy: 0.6827
31 [epoch 31] loss: 0.00481 accuracy: 0.7867 val loss: 0.00757 val accuracy: 0.6784
32 [epoch 32] loss: 0.00468 accuracy: 0.7914 val loss: 0.00775 val accuracy: 0.6796
33 [epoch 33] loss: 0.00456 accuracy: 0.7984 val loss: 0.00821 val accuracy: 0.6593
34 [epoch 34] loss: 0.00447 accuracy: 0.8022 val loss: 0.00774 val accuracy: 0.6772
35 [epoch 35] loss: 0.00435 accuracy: 0.8074 val loss: 0.00784 val accuracy: 0.6806
36 [epoch 36] loss: 0.00422 accuracy: 0.8147 val loss: 0.00810 val accuracy: 0.6797
37 [epoch 37] loss: 0.00412 accuracy: 0.8184 val loss: 0.00803 val accuracy: 0.6729
38 [epoch 38] loss: 0.00399 accuracy: 0.8241 val loss: 0.00788 val accuracy: 0.6894
39 [epoch 39] loss: 0.00388 accuracy: 0.8268 val loss: 0.00814 val accuracy: 0.6755
40 [epoch 40] loss: 0.00379 accuracy: 0.8336 val loss: 0.00815 val accuracy: 0.6795
41 [epoch 41] loss: 0.00367 accuracy: 0.8383 val loss: 0.00827 val accuracy: 0.6787
42 [epoch 42] loss: 0.00357 accuracy: 0.8418 val loss: 0.00879 val accuracy: 0.6636
43 [epoch 43] loss: 0.00346 accuracy: 0.8463 val loss: 0.00859 val accuracy: 0.6783
44 [epoch 44] loss: 0.00336 accuracy: 0.8512 val loss: 0.00883 val accuracy: 0.6674
45 [epoch 45] loss: 0.00331 accuracy: 0.8524 val loss: 0.00933 val accuracy: 0.6644
46 [epoch 46] loss: 0.00319 accuracy: 0.8580 val loss: 0.00900 val accuracy: 0.6700
47 [epoch 47] loss: 0.00307 accuracy: 0.8630 val loss: 0.00949 val accuracy: 0.6604
48 [epoch 48] loss: 0.00298 accuracy: 0.8682 val loss: 0.00949 val accuracy: 0.6679
49 [epoch 49] loss: 0.00291 accuracy: 0.8707 val loss: 0.00897 val accuracy: 0.6744
50 [epoch 50] loss: 0.00279 accuracy: 0.8781 val loss: 0.00953 val accuracy: 0.6761
```

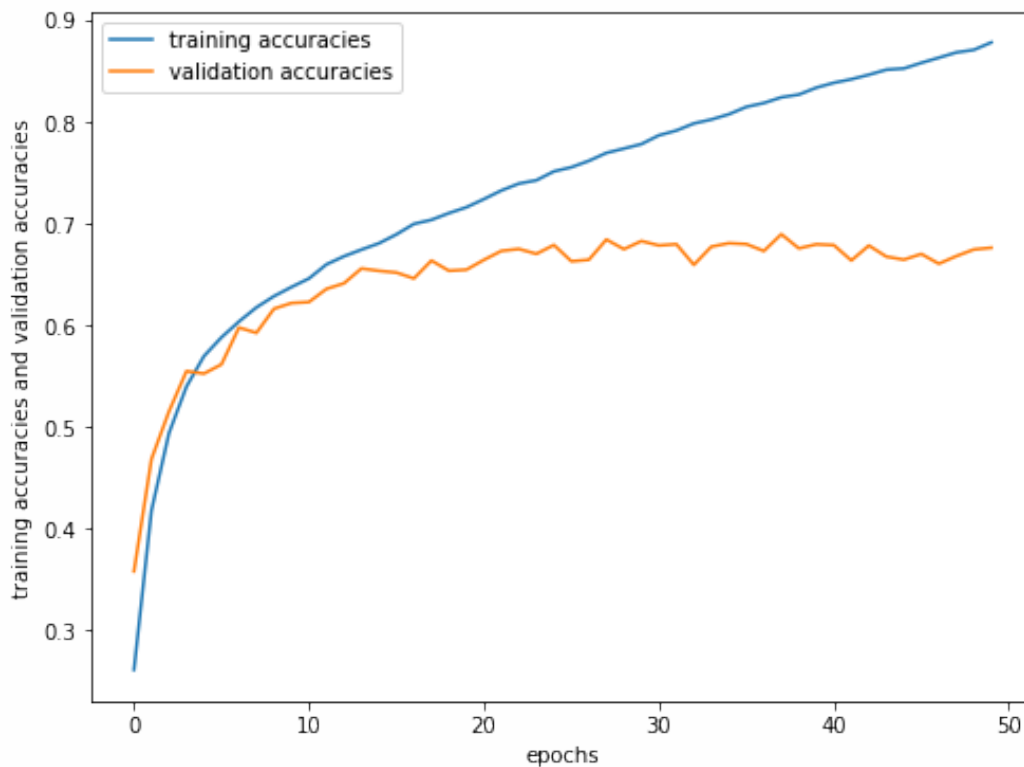
```
1 loss_list1 = [loss_list[i].item() for i in range(len(loss_list))]
2 loss_list_val1 = [loss_list_val[i].item() for i in range(len(loss_list_val))]
3 acc_list1 = [acc_list[i].item() for i in range(len(acc_list))]
4 acc_list_val1 = [acc_list_val[i].item() for i in range(len(acc_list_val))]
5
```

```
1 ##TODO: Plot the training losses and validation losses
2 X = range(max_epochs)
3 y1 = loss_list1
4 y2 = loss_list_val1
5 fig1, ax1 = plt.subplots(figsize=(8,6))
6 plt.plot(X, y1)
7 plt.plot(X, y2)
8 plt.xlabel('epochs')
9 plt.ylabel('training losses and validation losses')
10 plt.legend(['training losses', 'validation losses']);
11
```



```
1  ##TODO: Plot the training accuracies and validation accuracies
2  X = range(max_epochs)
3  y3 = acc_list1
4  y4 = acc_list_val1
5  fig2, ax2 = plt.subplots(figsize=(8,6))
6  plt.plot(X, y3)
7  plt.plot(X, y4)
8  plt.xlabel('epochs')
9  plt.ylabel('training accuracies and validation accuracies')
10 plt.legend(['training accuracies', 'validation accuracies']);
```





This model overfits on the data. Because the training accuracy is higher than validation accuracy.

```
1  #Test
2  true_labels = []
3  predictions = []
4  correct_test = 0
5  model_test = net.eval()
6  with torch.no_grad():
7      for batch_idx, (data, label) in enumerate(test_loader):
8          data, label = data.to(device), label.to(device)
9          ##TODO: pass the data into the network and store the output
10         output_test = model_test(data)
11
12         ##TODO: Get the prediction from the output
13         _, predicted_test = torch.max(output_test.data, 1)
14
15
16         ##TODO: Calculate the correct number and add the number to correct_test
17         correct_test += (predicted_test == label).sum()
18
19
20         ##TODO: update predictions list and true label list
21         true_labels.append(label.item())
22         predictions.append(predicted_test.item())
23
```

```

24         ##We can directly append the value because here batch_size=1
25
26
27
28 print('Accuracy on the 10000 test images: %.2f %%' % (100 * correct_test /
    len(test_set)))

```

```

1 Accuracy on the 10000 test images: 67.55 %

```

```

1 ##TODO: print the confusion matrix of test set
2 ##You can use sklearn.metrics.confusion_matrix
3 confusion_matrix(true_labels, predictions)

```

```

1 array([[671, 28, 61, 13, 53, 11, 20, 22, 86, 35],
2        [ 18, 821, 3, 6, 6, 1, 16, 7, 34, 88],
3        [ 63, 11, 520, 61, 127, 39, 92, 58, 23, 6],
4        [ 21, 13, 69, 441, 105, 100, 144, 64, 24, 19],
5        [ 10, 10, 52, 42, 669, 14, 89, 91, 17, 6],
6        [ 11, 4, 60, 198, 95, 433, 80, 101, 7, 11],
7        [ 4, 7, 23, 44, 41, 13, 853, 8, 6, 1],
8        [ 13, 10, 23, 34, 74, 31, 15, 779, 4, 17],
9        [ 57, 57, 17, 10, 18, 2, 7, 6, 796, 30],
10       [ 26, 115, 9, 10, 14, 1, 6, 18, 29, 772]])

```

Cat and dog get confused most. This makes sense because cat looks very similar to dogs.

```

1 |

```