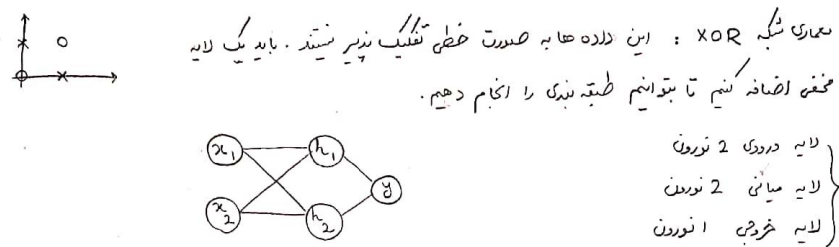
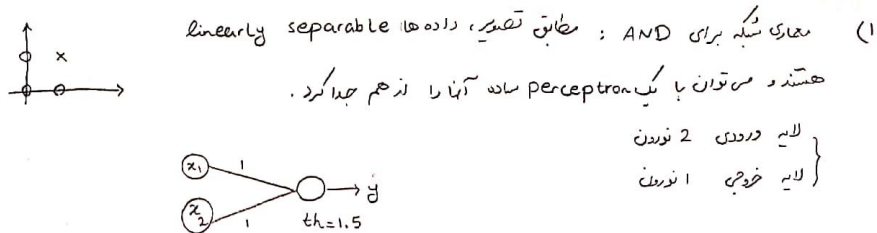


به نام خدا

تمرین چهارم یادگیری ماشین

غزل زمانی نژاد

1.



● الگوریتم backpropagation به دنبال بهینه کردن وزن های میان نودون ها است. در هر مرحله خطای پیش بینی نسبت به مقدار واقعی محاسبه می شود و خطا به عقب منتقل می شود تا وزن ها به سمت مقدار بهینه آپدیت شوند.

● در ابتدا وزن ها به صورت تصادفی مقداردهی می شوند. البته می توانیم از روش های مثل Xavier یا He هم استفاده کنیم که وزن ها را چیزی مقداردهی می کنند که از توزیع نرمال پیروی کند. بدین ترتیب از vanishing/exploding gradients تا حدی جلوگیری می شود.

● در یک epoch :

- (1) در forward pass، خروجی شبکه را بر اساس مقادیر فعلی وزن ها محاسبه می کنیم
 - (2) بر اساس loss function تعیین شده، میزان خطا بین خروجی و مقدار واقعی را بدست می آوریم.
 - (3) گرادیان خطا نسبت به وزن ها را حساب می کنیم و بر اساس فرمول $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$ وزن ها را آپدیت می کنیم.
- تا رسیدن به convergence، مراحل 1 تا 3 را تکرار می کنیم.

2. الف) استراتژی: توپ‌ها را در ۳ گروه تقسیم می‌کنیم. دو گروه را با ترازو مقایسه می‌کنیم: اگر وزن این دو گروه برابر باشد، توپ سنگین در گروه سوم است. در غیر این صورت توپ متفاوت در کفه سنگین‌تر ترازو است. سپس باید ۳ توپی که از گروه سنگین مانده را مقایسه کنیم که در اینجا هم مشابه قبل عمل می‌کنیم. دو توپ را روی ترازو می‌گذاریم: اگر وزن‌شان برابر باشد، توپ سوم توپ سنگین است و در غیر این صورت کفه سنگین‌تر نشان‌دهنده توپ متفاوت است. در هر بار مقایسه ۳ حالت ممکن است اتفاق بیفتد: وزن دو کفه برابر باشد، کفه راست سنگین‌تر باشد، کفه چپ سنگین‌تر باشد که می‌توانیم آن را کد k -nary با $k=3$ در نظر بگیریم. از آنجایی که ۲ آزمایش نیاز داریم پس در اینجا $i=2$ است پس

$$P(x=i) = 3^{-2} = 1/9$$

که برابر با احتمال خروج ۱ توپ از میان ۹ توپ است. سپس آنتروپی را محاسبه می‌کنیم:

$$-\sum_{i=1}^9 \frac{1}{9} \log_3 \frac{1}{9} = 2$$

پس حداقل آزمایش مورد نیاز ۲ است و با آنتروپی آن را اثبات کردیم.

ب) در آزمایش اول، از مقایسه هر دو گروه توپ سه‌تایی، به این میزان gain می‌رسیم:

$$H(X) - H(X|Y1) = 2 - 1 = 1$$

که این بیشترین مقدار است زمانی که خروجی سه حالت می‌تواند داشته باشد. در آزمایش دوم، مقایسه هر دو توپ از گروه سه‌تایی به میزان

$$H(X|Y1) - H(X|Y2, Y1) = 1 - 0 = 1$$

اطلاعات دارد. بنابراین ID3 به درخت تصمیم بهینه خواهد رسید.

$$3) \text{ total entropy} = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$$

Feature: Genre = {Drama, Horror}

$$\left. \begin{aligned} e(S_{\text{Drama}}) &= -\frac{3}{3} \log \frac{3}{3} = 0 \\ e(S_{\text{Horror}}) &= -1 \log 1 = 0 \end{aligned} \right\} \rightarrow \text{gain}(S, \text{Genre}) = 0.811 - 0 = \boxed{0.811}$$

Feature: Age Rating = {PG-13, G}

$$\left. \begin{aligned} e(S_{\text{PG-13}}) &= -\frac{3}{3} \log \frac{3}{3} = 0 \\ e(S_G) &= -1 \log 1 = 0 \end{aligned} \right\} \rightarrow \text{gain}(S, \text{Age}) = \boxed{0.811}$$

Feature: Language = {French, English}

$$\left. \begin{aligned} e(S_{\text{French}}) &= -1 \log 1 = 0 \\ e(S_{\text{English}}) &= -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \approx 0.923 \end{aligned} \right\} \rightarrow \text{gain}(S, \text{Language}) = 0.811 - \frac{2}{4} \cdot 0.923 = \boxed{0.118}$$

Feature: Source = {Book Adaptation}

$$e(S_{\text{Book}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \approx 0.811 \rightarrow \text{gain}(S, \text{Source}) = 0.811 - 0.811 = \boxed{0}$$

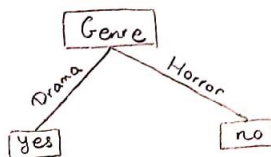
Feature: Film Location = {Canada, USA}

$$\left. \begin{aligned} e(S_{\text{Canada}}) &= -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \approx 0.923 \\ e(S_{\text{USA}}) &= -1 \log 1 = 0 \end{aligned} \right\} \rightarrow \text{gain}(S, \text{Location}) = 0.811 - \frac{2}{4} \times 0.923 = \boxed{0.118}$$

Feature: Studio = {Warner, A24}

$$\left. \begin{aligned} e(S_{\text{Warner}}) &= -\frac{2}{2} \log \frac{2}{2} = 0 \\ e(S_{\text{A24}}) &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1 \end{aligned} \right\} \rightarrow \text{gain}(S, \text{Studio}) = 0.811 - \frac{1}{2} \cdot 1 = \boxed{0.311}$$

در این مرحله بهترین gain را Genre دارد. پس انتخاب می‌کنیم Genre



$$\therefore \text{total entropy} = \frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.97$$

Feature: Genre = { Drama, Horror }

$$\left. \begin{aligned} e(S_{\text{Drama}}) &= -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.811 \\ e(S_{\text{Horror}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Genre}) = 0.97 - \frac{4}{5} \cdot 0.811 = \boxed{0.321}$$

Feature: Age Rating = { PG-13, G }

$$\left. \begin{aligned} e(S_{\text{PG-13}}) &= -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.811 \\ e(S_{\text{Horror}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Age}) = 0.97 - \frac{4}{5} \times 0.811 = \boxed{0.321}$$

Feature: Language = { French, English }

$$\left. \begin{aligned} e(S_{\text{French}}) &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1 \\ e(S_{\text{English}}) &= -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.923 \end{aligned} \right\} \text{gain}(S, \text{Language}) = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.923 = \boxed{0.016}$$

Feature: Source = { Book Adaptation, Original Screenplay }

$$\left. \begin{aligned} e(S_{\text{Book}}) &= -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.811 \\ e(S_{\text{screenplay}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Source}) = 0.97 - \frac{4}{5} \times 0.811 = \boxed{0.321}$$

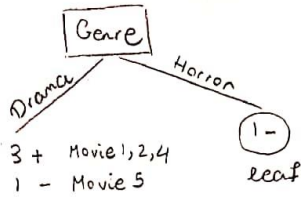
Feature: Film Location = { Canada, USA }

$$\left. \begin{aligned} e(S_{\text{Canada}}) &= -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} = 1 \\ e(S_{\text{USA}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Location}) = 0.97 - \frac{4}{5} \cdot 1 = \boxed{0.17}$$

Feature: Studio = { Warner, A24 }

$$\left. \begin{aligned} e(S_{\text{Warner}}) &= -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.923 \\ e(S_{\text{A24}}) &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1 \end{aligned} \right\} \text{gain}(S, \text{Studio}) = 0.97 - \frac{3}{5} \times 0.923 - \frac{2}{5} \times 1 = \boxed{0.016}$$

بیشترین gain 3 درخت : Genre, Source, Age, Genre : انتخاب مناسب



Movie 1, 2, 4, 5

- ⊕ 1 PG-13 French Book Canada Warner
- ⊕ 2 PG-13 English Book Canada Warner
- ⊕ 4 PG-13 English Book Canada A24
- ⊖ 5 PG-13 French Screenplay Canada Warner

$$\text{total entropy} = \frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \approx 0.811$$

Feature: Age = {PG-13}

$$e(S_{\text{PG-13}}) = \frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \approx 0.811 \rightarrow \text{gain}(S, \text{Age}) = 0.811 - 0.811 = 0$$

Feature: Language = {English, French}

$$\left. \begin{aligned} e(S_{\text{English}}) &= \frac{3}{3} \log \frac{3}{3} = 0 \\ e(S_{\text{French}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Language}) = 0.811$$

Feature: Source = {Book, Screenplay}

$$\left. \begin{aligned} e(S_{\text{Book}}) &= \frac{3}{3} \log \frac{3}{3} = 0 \\ e(S_{\text{screenplay}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Source}) = 0.811$$

Feature: Location = {Canada}

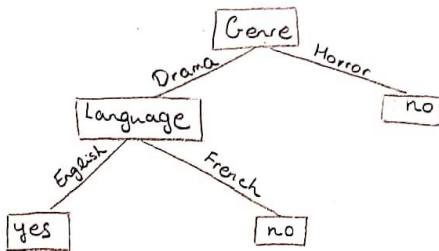
$$e(S_{\text{Canada}}) = \frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \approx 0.811 \rightarrow \text{gain}(S, \text{Location}) = 0$$

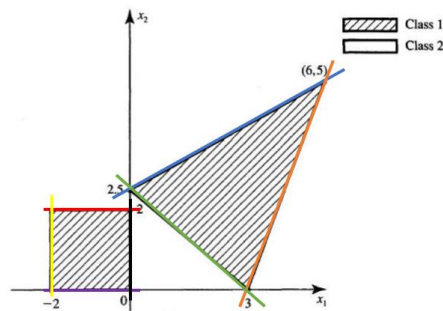
Feature: Studio = {Warner, A24}

$$\left. \begin{aligned} e(S_{\text{Warner}}) &= \frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \approx 0.923 \\ e(S_{\text{A24}}) &= -1 \log 1 = 0 \end{aligned} \right\} \text{gain}(S, \text{Studio}) = 0.811 - \frac{3}{4} \times 0.923 \approx 0.118$$

بسترن gain، و درستی Source، Language، و انتخاب می کنیم ← Language

درخت کلاسی





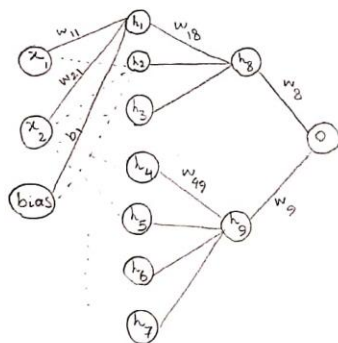
$$4) w_1 x_1 + w_2 x_2 + b = 0$$

معادله حرکت از خطوط را بدست می آوریم. پس آنها را ترتیب می دهیم.

$$\left. \begin{aligned} \text{خط آبی: } y &= \frac{2.5}{6} x_1 + 2.5 \rightarrow 2.5 x_1 - 6 x_2 + 15 \geq 0 \\ \text{خط سبز: } y &= \frac{5}{3} x_1 - 5 \rightarrow 5 x_1 - 3 x_2 - 15 \leq 0 \\ \text{خط سبز: } y &= \frac{-2.5}{3} x_1 + 2.5 \rightarrow 2.5 x_1 + 3 x_2 - 7.5 \leq 0 \end{aligned} \right\} \begin{array}{l} \text{استرایک} \\ \text{ناحیه 3} \end{array}$$

$$\left. \begin{aligned} \text{خط قرمز: } y &= 2 \rightarrow 0 x_1 + x_2 - 2 \leq 0 \\ \text{خط بنفش: } y &= 0 \rightarrow 0 x_1 + x_2 + 0 \geq 0 \\ \text{خط مشکی: } x &= 0 \rightarrow x_1 + 0 x_2 + 0 \leq 0 \\ \text{خط زرد: } x &= -2 \rightarrow x_1 + 0 x_2 + 2 \geq 0 \end{aligned} \right\} \begin{array}{l} \text{استرایک} \\ \text{ناحیه 4} \end{array}$$

لایه هیدین اول:



$$w_{11} = 2.5, w_{21} = -6, b_1 = 15.5$$

$$w_{12} = -5, w_{22} = 3, b_2 = 15.5$$

$$w_{13} = -2.5, w_{23} = -3, b_3 = 8$$

$$w_{14} = 0, w_{24} = -1, b_4 = 2.5$$

$$w_{15} = 0, w_{25} = 1, b_5 = 0.5$$

$$w_{16} = -1, w_{26} = 0, b_6 = -0.5$$

$$w_{17} = 1, w_{27} = 0, b_7 = 2.5$$

وزن های نورون خروجی:

$$w_8 = \frac{1}{2}, w_9 = \frac{1}{2}$$

active سرن قاسی نورون ها:

$$y = \begin{cases} 1 & w^T x + b \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

لایه هیدین دوم:

$$w_{18} = \frac{1}{6}, w_{19} = 0$$

$$w_{28} = \frac{1}{6}, w_{29} = 0$$

$$w_{38} = \frac{1}{6}, w_{39} = 0$$

$$w_{48} = 0, w_{49} = \frac{1}{8}$$

$$w_{58} = 0, w_{59} = \frac{1}{8}$$

$$w_{68} = 0, w_{69} = \frac{1}{8}$$

$$w_{78} = 0, w_{79} = \frac{1}{8}$$

۵.۱ .

۵.۱: بخشی از داده آموزش را به عنوان اعتبارسنجی استفاده می‌کنیم. همچنین از transform برای نرمالیزه کردن تصاویر استفاده می‌کنیم.

```
1 transform = transforms.Compose(  
2     [transforms.ToTensor(),  
3     transforms.Normalize((0.5,), (0.5,))] )  
4  
5 batch_size = 32  
6  
7 dataset = torchvision.datasets.MNIST(root='./data', train=True, download=True, transform=transform)  
8 train_set, val_set = torch.utils.data.random_split(dataset, [0.8, 0.2])  
9  
10 train_loader = torch.utils.data.DataLoader(train_set, batch_size=batch_size, shuffle=True)  
11 val_loader = torch.utils.data.DataLoader(val_set, batch_size=batch_size, shuffle=False)  
12  
13 test_set = torchvision.datasets.MNIST(root='./data', train=False, download=True, transform=transform)  
14 test_loader = torch.utils.data.DataLoader(test_set, batch_size=batch_size, shuffle=False)
```

پارامترهای شبکه و لایه های مختلف: از یک mlp دو لایه استفاده می‌کنیم. چون داده ورودی تصویر است و دو بعد دارد در ابتدا از یک لایه flatten استفاده می‌کنیم. سپس از یک لایه fully connected با 28×28 نورون ورودی و ۱۲۸ نورون خروجی، تابع فعالساز relu ، لایه fully connected دیگر با ۱۲۸ نورون ورودی و ۱۰ نورون خروجی استفاده می‌کنیم.

```
1 class MLP(nn.Module):  
2  
3     def __init__(self, input_dim, h1, output_dim):  
4         super().__init__()  
5         self.flatten = nn.Flatten()  
6         self.fc1 = nn.Linear(input_dim, h1)  
7         self.relu = nn.ReLU()  
8         self.fc2 = nn.Linear(h1, output_dim)  
9  
10    def forward(self, x):  
11        x = self.flatten(x)  
12        x = self.fc1(x)  
13        x = self.relu(x)  
14        x = self.fc2(x)  
15        return x
```

تابع train:

```

1 def train(train_loader, val_loader, model, loss_fn, optimizer, epochs=10):
2
3     train_losses = []
4     val_losses = []
5     train_accuracies = []
6     val_accuracies = []
7
8     for epoch in range(epochs):
9
10        train_loss = 0
11        train_acc = 0
12
13        # train mode
14        model.train()
15        # iterate through batches
16        for b, data in enumerate(train_loader):
17            inputs, labels = data
18            inputs = inputs.to(device)
19            labels = labels.to(device)
20            # zero the grads
21            optimizer.zero_grad()
22            # forward pass
23            outputs = model(inputs)
24            # find predictions for calculating accuracy
25            preds = torch.argmax(outputs, dim=1)
26            # calculate loss
27            loss = loss_fn(outputs, labels)
28            # backpropagation
29            loss.backward()
30            # update parameters
31            optimizer.step()
32
33            train_loss += loss.item()
34            train_acc += torch.sum(labels == preds).item()
35
36        train_loss /= len(train_loader)
37        train_acc /= len(train_loader.dataset)

```

```

38
39        # val mode
40        model.eval()
41        val_loss = 0
42        val_acc = 0
43        with torch.no_grad():
44            for b, data in enumerate(val_loader):
45                inputs, labels = data
46                inputs = inputs.to(device)
47                labels = labels.to(device)
48                outputs = model(inputs)
49                preds = torch.argmax(outputs, dim=1)
50                loss = loss_fn(outputs, labels)
51                val_loss += loss.item()
52                val_acc += torch.sum(labels == preds).item()
53
54        val_loss /= len(val_loader)
55        val_acc /= len(val_loader.dataset)
56
57        train_losses.append(train_loss)
58        val_losses.append(val_loss)
59        train_accuracies.append(train_acc)
60        val_accuracies.append(val_acc)
61
62        print(f'epoch {epoch}: train loss: {train_loss:.4f}, val loss: {val_loss:.4f}')
63
64    return train_losses, val_losses, train_accuracies, val_accuracies

```

مدل را در ۱۰ اپیاک و با نرخ آموزش ۰,۰۰۱ آموزش می‌دهیم. نتیجه آموزش:

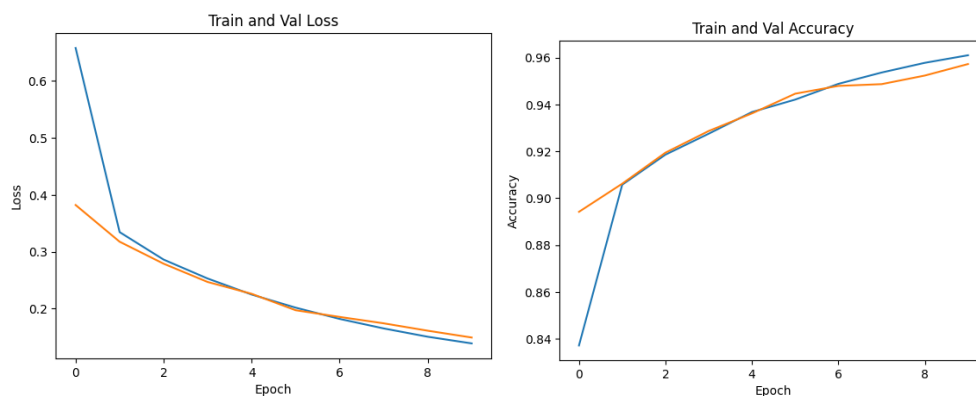
```

1 mlp = MLP(28*28, 128, 10).to(device)
2 loss_fn = nn.CrossEntropyLoss()
3 optimizer = torch.optim.Adam(mlp.parameters(), lr=1e-4)
4
5 train_losses, val_losses, train_accuracies, val_accuracies = train(train_loader, val_loader, mlp, loss_fn, optimizer, 10)

```

epoch 0: train loss: 0.6577, val loss: 0.3820
epoch 1: train loss: 0.3344, val loss: 0.3176
epoch 2: train loss: 0.2862, val loss: 0.2799
epoch 3: train loss: 0.2530, val loss: 0.2469
epoch 4: train loss: 0.2248, val loss: 0.2261
epoch 5: train loss: 0.2018, val loss: 0.1973
epoch 6: train loss: 0.1820, val loss: 0.1855
epoch 7: train loss: 0.1653, val loss: 0.1743
epoch 8: train loss: 0.1507, val loss: 0.1613
epoch 9: train loss: 0.1391, val loss: 0.1494

Learning curve



دقت طبقه‌بند: توانستیم به دقت ۹۵٪ بر روی داده آزمون برسیم.

```
1 from sklearn.metrics import confusion_matrix
2
3 def test(test_loader, model):
4     model.eval()
5
6     all_labels = []
7     all_predictions = []
8
9     with torch.no_grad():
10         for b, data in enumerate(test_loader):
11             inputs, labels = data
12             inputs = inputs.to(device)
13             labels = labels.to(device)
14             outputs = model(inputs)
15             preds = torch.argmax(outputs, dim=1)
16
17             all_labels.extend(labels.cpu().numpy())
18             all_predictions.extend(preds.cpu().numpy())
19
20     return all_labels, all_predictions
21
22 all_labels, all_predictions = test(test_loader, mlp)
23 test_acc = (np.array(all_labels) == np.array(all_predictions)).sum() / len(all_labels)
24
25 print("test accuracy:", test_acc)
26 print(confusion_matrix(all_labels, all_predictions))
27
28 test accuracy: 0.9577
29 [[ 962   0   1   1   0   4   8   2   2   0]
30 [   0 1119   2   2   0   1   5   2   4   0]
31 [   5   5 978  11   3   1   6  10  12   1]
32 [   1   2   4 976   0   4   2   8   7   6]
33 [   1   1   4   0 926   0  10   3   2 35]
34 [   7   2   1  24   1 825  12   1  12   7]
35 [   6   3   0   1   5   7 931   0   5   0]
36 [   2   8  12   6   1   1   0 982   1  15]
37 [   4   2   3  16   6   4  10   9 915   5]
38 [   4   6   1  12   9   1   2  10   1 963]]
```

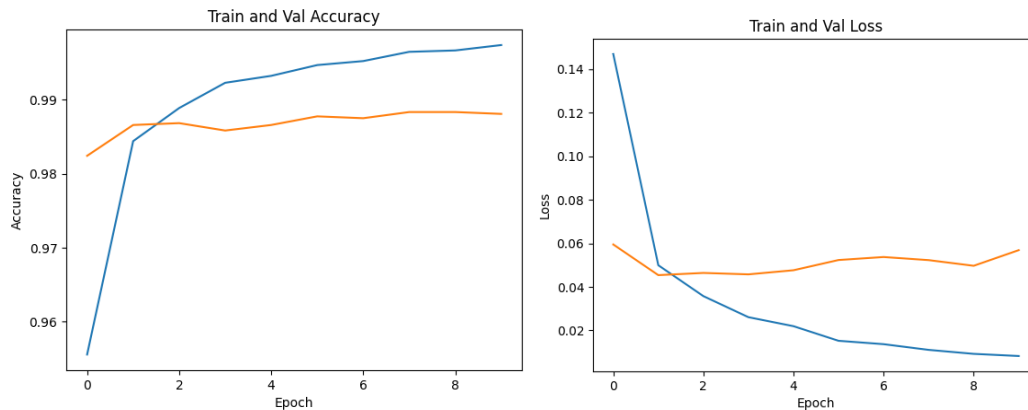
بهترین زمان متوقف کردن شبکه: مطابق نمودار خطا برای داده آموزش و اعتبارسنجی، از ایپاک پنجم به بعد میزان خطای آموزشی کاهش یافته ولی خطای اعتبارسنجی بسیار کم کاهش یافته. ممکن است مدل به سمت اورفیت شدن پیش رود پس بهتر است آموزش را در ایپاک ۵ متوقف کنیم.

۵,۲) پارامترهای شبکه و لایه های مختلف: شبکه cnn به صورت زیر است:

```
1 class CNN(nn.Module):
2     def __init__(self):
3         super(CNN, self).__init__()
4         self.conv1 = nn.Conv2d(1, 32, kernel_size=3, padding=1)
5         self.relu1 = nn.ReLU()
6         self.maxpool1 = nn.MaxPool2d(kernel_size=2, stride=2)
7
8         self.conv2 = nn.Conv2d(32, 64, kernel_size=3, padding=1)
9         self.relu2 = nn.ReLU()
10        self.maxpool2 = nn.MaxPool2d(kernel_size=2, stride=2)
11
12        self.flatten = nn.Flatten()
13        self.fc1 = nn.Linear(64 * 7 * 7, 10)
14
15    def forward(self, x):
16        x = self.conv1(x)
17        x = self.relu1(x)
18        x = self.maxpool1(x)
19
20        x = self.conv2(x)
21        x = self.relu2(x)
22        x = self.maxpool2(x)
23
24        x = self.flatten(x)
25        x = self.fc1(x)
26
27        return x
```

اولین لایه cnn کرنل 3×3 و ۳۲ کانال خروجی و تابع فعالسازی relu و بدنبال آن لایه max pooling است که به کاهش ابعاد کمک می کند. دومین لایه cnn کرنل 3×3 و ۶۴ کانال خروجی و تابع فعالسازی relu و بدنبال آن لایه max pooling است. در آخر از یک لایه تماماً متصل برای پیش بینی خروجی استفاده می کنیم که دارای ۱۰ نورون خروجی است.

Learning curve

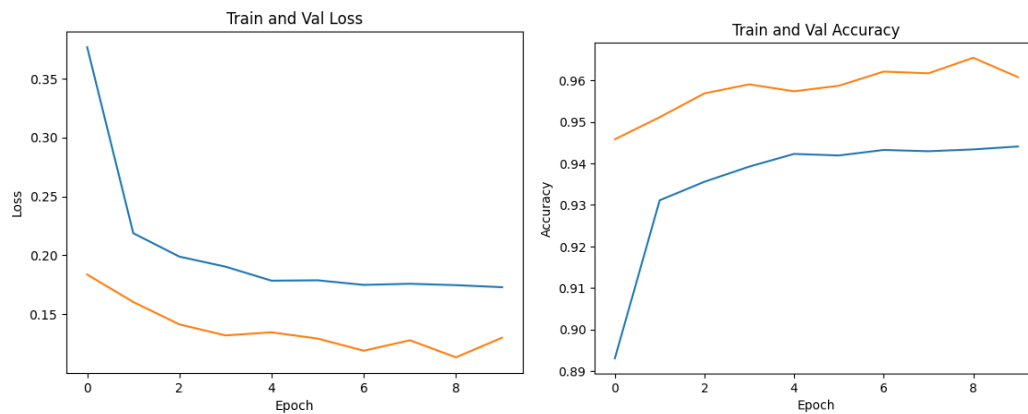


دقت طبقه بند:

```
test accuracy: 0.989
[[ 975   0   1   0   0   0   0   1   3   0]
 [   0 1130   2   0   0   0   1   1   1   0]
 [   0   0 1029   0   0   0   0   2   1   0]
 [   0   1   4 1001   0   4   0   0   0   0]
 [   0   0   1   0 975   0   0   0   1   5]
 [   0   1   0   5   0 882   3   0   1   0]
 [   3   1   1   0   2   2 946   0   3   0]
 [   0   4   20   0   0   0   0 1002   1   1]
 [   3   1   3   0   0   1   0   0   965   1]
 [   0   0   0   0  11   7   0   3   3 985]]
```

(۵,۳)

Learning curve



دقت طبقه‌بند:

test accuracy: 0.9624

```
[ [ 957  0  1  1  4  5  6  1  1  4]
 [  0 1133  1  0  1  0  0  0  0  0]
 [  5  9 935 24  9 21  5 13  5  6]
 [  0  0  1 991  1 13  0  1  0  3]
 [  0  5  0  0 967  0  4  0  1  5]
 [  1  0  0 30  1 853  1  4  1  1]
 [  8  8  1  1  7  8 923  0  1  1]
 [  0 19  3  5 12  1  0 973  0 15]
 [  6  1  2  5  6  6  8  1 923 16]
 [  6  4  0  5  9  9  0  5  2 969]]
```

6. در ابتدا داده را می‌خوانیم و فاصله‌های اضافه را از مقادیر دیتافریم حذف می‌کنیم.

```
1 import pandas as pd
2
3 columns = ['target', 'workclass', 'education', 'marital-status', 'occupation', 'relationship', 'race', 'sex', 'native-country']
4 train_data = pd.read_csv('adult.train.10k.discrete', header=None)
5 train_data.columns = columns
6
7 for col in columns:
8     train_data[col] = train_data[col].str.strip()
9 X_train = train_data.iloc[:, 1:]
10 y_train = train_data.iloc[:, 0]
11
12 test_data = pd.read_csv('adult.test.10k.discrete', header=None)
13 test_data.columns = columns
14
15 for col in columns:
16     test_data[col] = test_data[col].str.strip()
17 X_test = test_data.iloc[:, 1:]
18 y_test = test_data.iloc[:, 0]
```

یک کلاس نود داریم که برای ساختن گره با `property` های مقدار، ویژگی و فرزندان از آن استفاده می‌کنیم. از تابع `calculate entropy` برای محاسبه میزان آنتروپی بر حسب فرمول استفاده میکنیم. به کمک تابع `calculate gain`، میزان آنتروپی کل را بدست می‌آوریم و بر اساس ویژگی خواسته شده در ورودی، آنتروپی مربوط به آن ویژگی و داده‌ها را حساب می‌کنیم. در نهایت میزان اطلاعات را بازمی‌گردانیم.

```
1 import numpy as np
2
3 class Node:
4     def __init__(self, value=None, feature=None, children=None):
5         self.value = value
6         self.feature = feature
7         self.children = children
8
9
10 def calculate_entropy(X):
11     # H(S) = - sigma(p*log(p))
12     e = 0
13     labels = X.unique()
14     for l in labels:
15         p = len(X[X==l]) / len(X)
16         e -= p * np.log2(p) if p > 0 else 0
17     return e
18
19
20 def calculate_gain(X, y, feature):
21     total_entropy = calculate_entropy(y)
22     conditional_entropy = 0
23
24     unique_values = X[feature].unique()
25     for value in unique_values:
26         # find all rows with the specific value
27         subset = X[X[feature]==value]
28         subset_labels = y[X[feature]==value]
29         if not subset.empty:
30             e = calculate_entropy(subset_labels)
31             conditional_entropy += e * len(subset)/len(X)
32
33     return total_entropy - conditional_entropy
34
```

به کمک تابع `find_most_informative`، از میان ویژگی‌هایی که تا کنون استفاده نشده‌اند، ویژگی که بیشترین میزان اطلاعات را به ما می‌دهد پیدا می‌کنیم. از تابع `most_common_class` برای پیدا کردن کلاسی که بیشترین تکرار را دارد استفاده می‌شود. اصلی‌ترین تابع، `build_tree` است که به صورت `recursive` عمل می‌کند. شرط خاتمه (رسیدن به برگ) این است که در داده‌های موجود تمامی برچسب‌ها یکسان باشند و یا ویژگی‌ای باقی نمانده باشد. در غیر این صورت، بهترین ویژگی را پیدا می‌کنیم. و بر اساس مقادیر آن، باید ادامه درخت را تشکیل دهیم که به صورت بازگشتی این کار انجام می‌شود.

```

36 def find_most_informative(X, y, features):
37     most_gain = -float("inf")
38     best_feature = None
39
40     for feature in features:
41         g = calculate_gain(X, y, feature)
42         if g > most_gain:
43             best_feature = feature
44             most_gain = g
45
46     return best_feature
47
48 def most_common_class(y):
49     classes, counts = np.unique(y, return_counts=True)
50     idx = np.argmax(counts)
51     return classes[idx]
52
53 def build_tree(X, y, features_left):
54
55     # base case
56     if len(y.unique()) == 1:
57         return Node(value=y.iloc[0])
58
59     if len(features_left) == 0 or len(X) == 0:
60         return Node(value=most_common_class(y))
61
62     best_feature = find_most_informative(X, y, features_left)
63     unique_values = X[best_feature].unique()
64
65     children = {}
66     for val in unique_values:
67         # find all rows which their feature value is 'val'
68         subset = X[X[best_feature]==val]
69         labels = y[X[best_feature]==val]
70         if not subset.empty:
71             new_features = [f for f in features_left if f != best_feature]
72             children[val] = build_tree(subset, labels, new_features)
73
74     return Node(feature=best_feature, children=children)
75

```

تابع predict برای پیش‌بینی کردن داده جدید به کار می‌رود و با traverse کردن درخت تا رسیدن به برگ عمل می‌کند.

```

76 def predict(tree, data):
77     try:
78         if tree.value:
79             return tree.value
80
81         feature = tree.feature
82         value = data[feature]
83         return predict(tree.children[value], data)
84
85     except:
86         return '<=50K'

```

میزان دقت مدل برای داده آموزش:

```

1 from sklearn.metrics import classification_report, accuracy_score
2
3 preds_train = []
4 for i in range(len(X_train)):
5     preds_train.append(predict(root, X_train.iloc[i]))
6
7 print(accuracy_score(y_train, preds_train))
8 print(classification_report(y_train, preds_train))

```

	precision	recall	f1-score	support
<=50K	0.90	0.94	0.92	7550
>50K	0.79	0.67	0.73	2450
accuracy			0.88	10000
macro avg	0.84	0.81	0.82	10000
weighted avg	0.87	0.88	0.87	10000

میزان دقت مدل برای داده آزمون:

```
1 preds_test = []
2 for i in range(len(X_test)):
3     preds_test.append(predict(root, X_test.iloc[i]))
4
5 print(accuracy_score(y_test, preds_test))
6 print(classification_report(y_test, preds_test))
```

0.8097

	precision	recall	f1-score	support
<=50K	0.85	0.91	0.88	7539
>50K	0.65	0.50	0.56	2461
accuracy			0.81	10000
macro avg	0.75	0.71	0.72	10000
weighted avg	0.80	0.81	0.80	10000

به نظر می‌رسد در داده تست برای برخی ویژگی‌ها مقادیری وجود داشت که در داده آموزش نیامده بود.