

دانشگاه صنعتی خواجه نصیرالدین طوسی  
دانشکده مهندسی برق - گروه مهندسی کنترل

## درس یادگیری ماشین پاسخ مینی پروژه اول

ایمان گندمی	نام و نام خانوادگی
۴۰۲۲۳۷۰۴	شماره دانشجویی
فروردین ۱۴۰۳	تاریخ



## فهرست مطالب

۸	۱ سوال اول
۸	۱.۱ فرض کنید در یک مسئله طبقه بندی دو کلاسه، دو لایه انتهایی شبکه شما فعال ساز ReLU و سیگموید است. چه اتفاقی می افتد؟
۹	۲.۱ یک جایگزین برای ReLU در معادله ۱ آورده شده است. ضمن محاسبه گرادیان آن، حداقل یک مزیت آن نسبت به ReLU را توضیح دهید
۹	۳.۱ به کمک یک نورون ساده، پرسپترون و یا نورون Pitts-McCulloch شبکه ای طراحی کنید که ...
۲۰	۲ سوال دوم
۲۰	۱.۲ دیتاست CWRU Bearing که در «مینی پروژه شماره یک» با آن آشنا شدید را به خاطر آورید. علاوه بر دو کلاسی که در آن مینی پروژه در نظر گرفتید، با مراجعه به صفحه داده های عیب در حالت ۱۲، k دو کلاس دیگر نیز از ۶OR@X-۰۰۷ اضافه کنید. با انجام این کار یک کلاس داده سالم و سه کلاس طریق فایل های ۱۰۷B-X از داده های دارای سه عیب متفاوت خواهید داشت. در مورد این که هر فایل مربوط به چه نوع عیبی است به صورت کوتاه توضیح دهید. سپس در ادامه، تمام کارهایی که در بخش «۲» سوال دوم «مینی پروژه یک» برای استخراج ویژگی و آماده سازی دیتا انجام داده بودید را روی دیتاست جدید خود پیاده سازی کنید. در قسمت تقسیم بندی داده ها، یک بخش برای «اعتبارسنجی» به بخش های «آموزش» و «آزمون» اضافه کنید و توضیح دهید که کاربرد این بخش چیست.
۲۰	۲.۲ یک مدل MLP ساده با ۲ لایه پنهان یا بیش تر بسازید. بخشی از داده های آموزش را برای اعتبارسنجی کنار بگذارید و با انتخاب بهینه ساز و تابع اتلاف مناسب، مدل را آموزش دهید. نمودارهای اتلاف و Accuracy مربوط به آموزش و اعتبارسنجی را رسم و نتیجه را تحلیل کنید. نتیجه تست مدل روی داد های آزمون را با استفاده ماتریس درهم ریختگی و report-classification نشان داده و نتایج به صورت دقیق تحلیل کنید.
۲۷	۳.۲ فرآیند سوال قبل را با یک بهینه ساز و تابع اتلاف جدید انجام داده و نتایج را مقایسه و تحلیل کنید. بررسی کنید که آیا تغییر تابع اتلاف می تواند در نتیجه اثرگذار باشد؟
۳۲	۴.۲ در مورد Cross-validation K-Fold Stratified و Cross-validation K-Fold و مزایای هر یک توضیح دهید. سپس با ذکر دلیل، یکی از این روش ها را انتخاب کرده و بخش «۲» سوال سوم را با آن پیاده سازی کنید و نتایج خود را تحلیل کنید.
۳۴	۳ سوال سوم
۳۸	۱.۳ با استفاده از ماتریس درهم ریختگی و حداقل سه شاخصه ارزیابی مربوط به وظیفه طبقه بندی، عمل کرد درخت آموزش داده شده خود را روی بخش آزمون داده ها ارزیابی کنید و نتایج را به صورت دقیق گزارش کنید. تأثیر مقادیر کوچک و بزرگ حداقل دو فرایامتر را بررسی کنید. تغییر فرایامترهای مربوط به هرس کردن چه تأثیری روی نتایج دارد و مزیت آن چیست؟
۴۳	۲.۳ توضیح دهید که روش هایی مانند جنگل تصادفی و AdaBoost چگونه می توانند به بهبود نتایج کمک کنند. سپس، با انتخاب یکی از این روش ها و استفاده از فرایامترهای مناسب، سعی کنید نتایج پیاده سازی در مراحل قبلی را ارتقاء دهید.
۴۶	





## فهرست تصاویر

۹	.....	توابع فعال ساز Relu و ELU	۱
۱۵	.....	بلوک دیاگرام مدل طبقه بند خطی	۲
۱۶	.....	بلوک دیاگرام مدل طبقه بند خطی	۳
۱۶	.....	بلوک دیاگرام مدل طبقه بند خطی	۴
۱۷	.....	بلوک دیاگرام مدل طبقه بند خطی	۵
۲۰	.....	خروجی هر سه حالت استفاده از دو تابع فعال ساز مختلف	۶
۲۲	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۷
۲۲	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۸
۲۳	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۹
۲۶	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۰
۲۶	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۱
۲۷	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۲
۲۹	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۳
۳۰	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۴
۳۱	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۵
۳۱	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۶
۳۲	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۷
۳۳	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۸
۳۳	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۹
۳۴	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۰
۳۷	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۱
۳۸	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۲
۳۹	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۳
۴۰	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۴
۴۱	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۵
۴۲	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۶
۴۳	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۷
۴۴	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۸
۴۵	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۹
۴۶	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۰
۴۸	.....	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۱
۴۸	.....	بلوک دیاگرام مدل طبقه بند خطی	۳۲
۴۹	.....	بلوک دیاگرام مدل طبقه بند خطی	۳۳
۵۰	.....	بلوک دیاگرام مدل طبقه بند خطی	۳۴



۵۲	.....	بلوک دیاگرام مدل طبقه‌بند خطی	۳۵
۵۵	.....	بلوک دیاگرام مدل طبقه‌بند خطی	۳۶
۵۶	.....	بلوک دیاگرام مدل طبقه‌بند خطی	۳۷



## فهرست جداول



## فهرست برنامه‌ها

۹	..... (Python) libraries import	۱
۱۱	..... (Python) libraries import	۲
۱۱	..... (Python) libraries import	۳
۱۲	..... (Python) libraries import	۴
۱۳	..... (Python) libraries import	۵
۱۳	..... (Python) libraries import	۶
۱۴	..... (Python) libraries import	۷
۱۷	..... (Python) libraries import	۸
۱۸	..... (Python) libraries import	۹
۱۹	..... (Python) libraries import	۱۰
۲۱	..... library and dataset import	۱۱
۲۲	..... library and dataset import	۱۲
۲۲	..... library and dataset import	۱۳
۲۳	..... library and dataset import	۱۴
۲۴	..... library and dataset import	۱۵
۲۵	..... library and dataset import	۱۶
۲۵	..... library and dataset import	۱۷
۲۶	..... library and dataset import	۱۸
۲۶	..... library and dataset import	۱۹
۲۷	..... library and dataset import	۲۰
۲۸	..... library and dataset import	۲۱
۲۸	..... library and dataset import	۲۲
۲۹	..... library and dataset import	۲۳
۲۹	..... library and dataset import	۲۴
۳۰	..... library and dataset import	۲۵
۳۵	..... library and dataset import	۲۶
۳۵	..... library and dataset import	۲۷
۳۸	..... library and dataset import	۲۸
۳۸	..... library and dataset import	۲۹
۴۰	..... library and dataset import	۳۰
۴۱	..... library and dataset import	۳۱
۴۱	..... library and dataset import	۳۲
۴۲	..... library and dataset import	۳۳
۴۲	..... library and dataset import	۳۴



۴۴	.....	library and dataset import	۳۵
۴۷	.....	library and dataset import	۳۶
۴۸	.....	(Python) libraries import	۳۷
۴۸	.....	(Python) libraries import	۳۸
۴۹	.....	(Python) libraries import	۳۹
۴۹	.....	(Python) libraries import	۴۰
۵۰	.....	(Python) libraries import	۴۱
۵۱	.....	(Python) libraries import	۴۲
۵۱	.....	(Python) libraries import	۴۳
۵۲	.....	(Python) libraries import	۴۴
۵۲	.....	(Python) libraries import	۴۵
۵۳	.....	(Python) libraries import	۴۶
۵۴	.....	(Python) libraries import	۴۷
۵۶	.....	(Python) libraries import	۴۸





لینک پوشه گیت هاب:

[https://github.com/ImanGandomi/MachineLearning2024/tree/main/ml\\_MiniProject\\_1](https://github.com/ImanGandomi/MachineLearning2024/tree/main/ml_MiniProject_1)

لینک گوگل کولب:

[https://colab.research.google.com/drive/1AjG5WVYKpsJYs8rkoI94V5ZuJdnVG3\\_m?usp=sharing](https://colab.research.google.com/drive/1AjG5WVYKpsJYs8rkoI94V5ZuJdnVG3_m?usp=sharing)

## ۱ سوال اول

۱.۱ فرض کنید در یک مسئله طبقه بندی دوکلاسه، دو لایه انتهایی شبکه شما فعال ساز ReLU و سیگموید است. چه اتفاقی می افتد؟

ابتدا در مورد مزیت استفاده از هر یک از این توابع فعال ساز می توان صحبت کرد. تابع فعال ساز ReLU به دلیل سادگی و اثرگذاری در یادگیری عمیق شناخته می شود. این تابع باعث ایجاد حالت غیرخطی می شود به این صورت که اگر مقدار ورودی مثبت باشد آن را عبور داده و در غیر این صورت صفر خروجی می دهد. این تابع باعث غلبه بر مشکل محوشدگی گرادیان می شود. این تابع نسبت به سایر توابع مثل سیگموید به مرحله backpropagation بهتر کمک کرده و برای ورودی های بزرگتر خروجی ها را کوچک نمی کند و در نتیجه محوشدگی گرادیان کمتر بروز می کند. کاهش بار محاسباتی به دلیل سادگی این گرادیان از دیگر مزایای استفاده از آن است. همچنین این سادگی باعث می شد که میزان سرعت همگرایی شبکه به شدت افزایش یابد.

تابع فعال ساز sigmoid باعث هدایت خروجی به بین مقادیر ۰ و ۱ می شود. خروجی این تابع به عنوان احتمال تعلق به کلاس خاص تعبیر می شود و معمولاً در لایه آخر برای انجام طبقه بندی نهایی مورد استفاده قرار می گیرد. استفاده از این تابع باعث ایجاد حالت smooth gradient می شود و تا حدودی باعث آسانی در مرحله backpropagation می شود. البته استفاده از این تابع به دلیل کوچک کردن ورودی های خود به مقدار بین ۰ و ۱ باعث ایجاد محوشدگی گرادیان می شود.

در مورد سوال مطرح شده می توان اینگونه پاسخ داد که استفاده از این نوع ساختار در لایه های انتهایی مدل طبقه بند به نوعی رایج است. مزایای استفاده از این ساختار را با توجه به توضیحات گفته شده می توان این طور مطرح کرد: برای لایه ماقبل آخر با تابع فعال ساز ReLU ایجاد غیرخطی گری و استخراج ویژگی: پیش از لایه آخر که کار طبقه بندی صورت می گیرد، استفاده از این تابع باعث می شود ویژگی هایی که خاصیت غیرخطی دارند به طور مناسب تری از یکدیگر قابلیت تفکیک پیدا کنند. همچنین پیش از ورودی به لایه طبقه بند، مدل با استفاده از این تابع فعال ساز قابلیت یادگیری ویژگی های سطح بالاتر را به دست می آورد

گرادیان و سرعت بالا: همانطور هم که قبلاً ذکر شد این تابع باعث عدم ایجاد محوشدگی گرادیان شده و آموزش شبکه را هموار می کند. همچنین به دلیل ساختار ساده خود باعث همگرایی سریع تر مدل می شود

برای لایه آخر با تابع فعال ساز: sigmoid

تفسیرپذیری: با استفاده از این تابع در لایه آخر که خروجی بین ۰ و ۱ ایجاد می کند، می توان از این خروجی به عنوان احتمال تعلق آن داده به کلاس خاص یاد کرد.

انتقال هموار: ترکیب این دو تابع فعال ساز باعث ایجاد انتقال هموار بین لایه های استخراج ویژگی و لایه طبقه بندی می شود. این شرایط تضمین می کند که ویژگی های آموخته شده مناسب هستند و منجر به پیش بینی مناسب می شود.

یک مورد که لازم است که به آن توجه کرد این است که با این که این دو تابع فعال ساز به طور معمول در شبکه های مختلف مورد استفاده قرار می گیرند، اما اگر به ساختار و خروجی هر یک دقت کنیم متوجه می شویم که پس از عبور داده ها از تابع فعال ساز ReLU مقادیر خروجی به ازای مقادیر مثبت، همان مقدار باقی می ماند در صورتی که مقادیر منفی در خروجی صفر می شود. حال این مقادیر



مثبت و یا صفر در تابع فعال‌ساز بعدی (سیگموید) که به ازای مقادیر مثبت مقداری بین ۰ و ۵۰٪ و خروجی می‌دهد و به ازای مقادیر منفی مقداری بین ۰ و ۵۰٪ خروجی می‌دهد، بین ۰ و ۱۰۰٪ می‌شوند. این فرآیند ممکن است کمی تصمیم‌گیری را لایه آخر محدود کند.

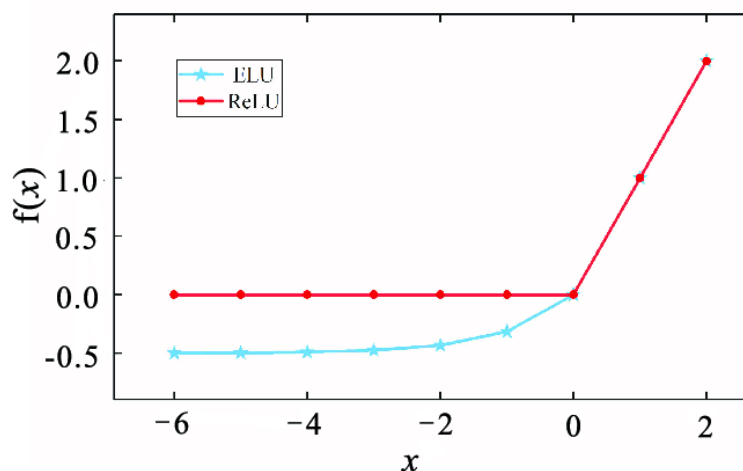
## ۲.۱ یک جایگزین برای ReLU در معادله ۱ آورده شده است. ضمن محاسبه گرادیان آن، حداقل یک مزیت آن نسبت به ReLU را توضیح دهید

گرادیان این تابع به صورت زیر است.

پارامتر  $\alpha$  در این تابع تعیین‌کننده اندازه خروجی تابع در مقادیر منفی ورودی است.  $\alpha$  معمولاً مقداری مثبت و کوچک در نظر گرفته می‌شود. مثلاً ۱.۰ یا ۰.۵۰.

یک مشکل مهم تابع ReLU مقادیر خروجی این تابع به ازای ورودی‌های منفی است. زمانی که ورودی منفی باشد، تابع ReLU مقدار خروجی نوروں را صفر می‌کند و برای همه به‌روزرسانی‌های بعدی همین مشکل باقی می‌ماند. این مشکل با عنوان dying ReLU مطرح است. تابع ELU با اجازه‌دادن عبور مقادیر منفی برای ورودی، این مشکل را برطرف می‌کند. این کار باعث فعال باقی ماندن نوروں می‌شود.

در شکل ۱ تصویری از این دو تابع و تفاوت در خروجی مقادیر منفی آورده شده است.



شکل ۱: توابع فعال‌ساز ReLU و ELU

## ۳.۱ به کمک یک نوروں ساده، پرسپترون و یا نوروں Pitts-McCulloch شبکه ای طراحی کنید که ...

ابتدا کتابخانه‌های مورد نیاز را فراخوانی می‌کنیم.

```
1 #import library
2 import numpy as np
3 import itertools
4 import matplotlib.pyplot as plt
```

Code 1: import libraries (Python)



سه نقطه داده شده برای سه راس مثلث مورد نظر به این صورت هستند.

$$(2, 2), \quad (1, 0), \quad (3, 0)$$

سه خط برای سه ضلع مثلث باید محاسبه شوند که هر یک از این سه خط را به صورت جداگانه به عنوان یک خط جداکننده دوکلاس می‌توان در نظر گرفت. برای هر یک از این سه طبقه‌بندی، یک نورون با دو ورودی و یک آستانه باید در نظر گرفت.

۱. برای خط گذرنده از دو نقطه  $(2, 2)$  و  $(1, 0)$ :

$$\begin{aligned} m &= \frac{y_2 - y_1}{x_2 - x_1} \\ &= \frac{0 - 2}{1 - 2} \\ &= \frac{-2}{-1} \\ &= 2 \end{aligned}$$

Using point-slope form with  $(2, 2)$ :

$$y - 2 = 2(x - 2)$$

$$y - 2 = 2x - 4$$

$$y = 2x - 2$$

۲. برای خط گذرنده از دو نقطه  $(2, 2)$  و  $(3, 0)$ :

$$\begin{aligned} m &= \frac{y_2 - y_1}{x_2 - x_1} \\ &= \frac{0 - 2}{3 - 2} \\ &= \frac{-2}{1} \\ &= -2 \end{aligned}$$

Using point-slope form with  $(2, 2)$ :

$$y - 2 = -2(x - 2)$$

$$y - 2 = -2x + 4$$

$$y = -2x + 6$$

۳. برای خط گذرنده از دو نقطه  $(1, 0)$  و  $(3, 0)$ :

$$\begin{aligned} m &= \frac{y_2 - y_1}{x_2 - x_1} \\ &= \frac{0 - 0}{3 - 1} \\ &= \frac{0}{2} \\ &= 0 \end{aligned}$$

Using point-slope form with  $(1, 0)$ :

$$y - 0 = 0(x - 1)$$

$$y = 0$$

در گام بعدی نرون McClulloch-Pitts را تعریف می‌کنیم. برای این کار یک کلاس تعریف می‌کنیم. در قسمت init این کلاس وزن‌ها و آستانه مورد نظر در مرحله مقداردهی اولیه داده می‌شود. در قسمت model مقادیر وزن و ورودی در هم ضرب شده و پس از مقایسه با مقدار آستانه به صورت دو کلاس ۰ و ۱ طبقه‌بندی می‌شوند.

```

1 #define muculloch pitts
2 class McCulloch_Pitts_neuron():
3     def __init__(self , weights , threshold):
4         self.weights = weights      #define weights
5         self.threshold = threshold   #define threshold
6
7     def model(self , x):
8         #define model with threshold
9         y = self.weights @ x
10        if y >= self.threshold:
11            return 1
12        else:
13            return 0

```

Code 2: import libraries (Python)

کد آمده در ۱۵ برای فراخوانی کلاس نرون McClulloch-Pitts با وزن‌ها و آستانه‌ها متناسب با هر کد جداکننده نوشته شده است. در این تابع ابتدا کلاس نرون فراخوانی شده و سپس ورودی‌های مورد نظر به کلاس داده می‌شود برای اینکه کار طبقه‌بندی صورت گیرد

```

1 #define model for dataset
2 def Area(x, y):
3     neur1 = McCulloch_Pitts_neuron([0, 1], 0)
4     neur2 = McCulloch_Pitts_neuron([2, -1], 2)
5     neur3 = McCulloch_Pitts_neuron([-2, -1], -6)
6     neur4 = McCulloch_Pitts_neuron([1, 1, 1], 3)
7
8     z1 = neur1.model(np.array([x, y]))
9     z2 = neur2.model(np.array([x, y]))

```



```
10 z3 = neur3.model(np.array([x, y]))
11 z4 = neur4.model(np.array([z1, z2, z3]))
12
13 # return str(z1) + str(z2)
14 return list([z1])
```

Code 3: import libraries (Python)

حالا ابتدا باید به هر یک از این سه خط به عنوان یک خط جداکننده دو کلاس نگاه کرد. به طوری که هر بار با نورون McCluloch-Pitts نقاط به دو سمت این خط جدا شوند. بنابراین سه مرتبه باید این طبقه‌بندی دو کلاسه را ابتدا انجام داد. سه خط محاسبه شده مقادیر وزن‌ها و آستانه مورد نظر برای طبقه‌بندی‌های مورد نظر را تشکیل می‌دهند. بنابراین در مرحله اول سه طبقه‌بندی با نورون McCluloch-Pitts با این مقادیر اولیه داریم:

$([0, 1], 0)$   
 $([2, -1], 2)$   
 $([-2, -1], -6)$

در مرحله بعد باید نتیجه این سه خط را با هم ترکیب کرد. این کار توسط یک نورون با سه ورودی صورت می‌گیرد. مقادیر اولیه این نورون به صورت زیر در نظر گرفته می‌شود:

$([1, 1, 1], 3)$

بنابراین با این مراحل ناحیه مورد نظر باید به صورت زیر حاصل شود.

```
1 # Initialize lists to store data points for different z5 values
2 red_points = []
3 green_points = []
4
5 # Evaluate data points using the Area function
6 for i in range(num_points):
7     z5_value = Area(x_values[i], y_values[i])
8     if z5_value == [0]: # z5 value is 0
9         red_points.append((x_values[i], y_values[i]))
10    else: # z5 value is 1
11        green_points.append((x_values[i], y_values[i]))
12
13 # Separate x and y values for red and green points
14 red_x, red_y = zip(*red_points)
```



```
15 green_x, green_y = zip(*green_points)
```

Code 4: import libraries (Python)

ابتدا برای خط

$$y = 0$$

در کد ۱۵ در خط اول قسمت فراخوانی، مقادیر  $[0, 1]$  را برای وزن‌دهی و مقدار  $0$  را برای آستانه مقداردهی می‌کنیم. مقدار  $z_1$  را به خروجی کلاس McClulloch-Pitts برای خط اول پس از دادن ورودی‌های مورد نظر نسبت می‌دهیم. در انتهای این تابع مقدار  $z_1$  را که خروجی تابع برای خط اول بود را برمی‌گردانیم. کد ۵ تعداد ۲۰۰۰ نقطه رندوم در دو بعد  $x$  و  $y$  تولید می‌کند.

```
1 # Generate random data points
2 num_points = 2000
3 x_values = np.random.uniform(0, 4, num_points)
4 y_values = np.random.uniform(-1, 3, num_points)
```

Code 5: import libraries (Python)

در قسمت بعدی باید تمامی نقاط تولید شده را به نوروں داده تا کار طبقه‌بندی برای این نقاط صورت بگیرد. در حلقه در نظر گرفته شده به تعداد نقاط رندوم کار طبقه‌بندی صورت می‌گیرد که اگر برای هر نقطه خروجی نوروں ۱ یا ۰ بود، آن نقطه به هر گروه سبز یا قرمز تعلق می‌یابد.

```
1 # Initialize lists to store data points for different z5 values
2 red_points = []
3 green_points = []
4
5 # Evaluate data points using the Area function
6 for i in range(num_points):
7     z5_value = Area(x_values[i], y_values[i])
8     if z5_value == [0]: # z5 value is 0
9         red_points.append((x_values[i], y_values[i]))
10    else: # z5 value is 1
11        green_points.append((x_values[i], y_values[i]))
12
13 # Separate x and y values for red and green points
14 red_x, red_y = zip(*red_points)
```



```
15 green_x, green_y = zip(*green_points)
```

Code 6: import libraries (Python)

سپس با استفاده از کد ۳۱ نقاط پیش‌بینی شده متعلق به هر دو کلاس قرمز و سبز و همچنین مثلث مورد نظر سوال (برای مقایسه درستی خروجی طبقه‌بندی و ناحیه مثلث) رسم شده است

```
1 # Plotting
2 plt.figure(figsize=(8, 6))
3 plt.scatter(red_x, red_y, color='red', label='z = 0')
4 plt.scatter(green_x, green_y, color='green', label='z = 1')
5 plt.xlabel('X values')
6 plt.ylabel('Y values')
7 plt.title('McCulloch-Pitts Neuron Outputs')
8
9
10 y_points = [0, 2, 0]
11 x_points = [1, 2, 3]
12
13 plt.plot(x_points[0:2], y_points[0:2], "black")
14 plt.plot(x_points[1:3], y_points[1:3], "black")
15 plt.plot([x_points[0], x_points[-1]], [y_points[0], y_points[-1]], "black",
16         label='Triangle')
17
18 plt.fill(x_points, y_points, hatch='/', edgecolor='pink', alpha=0.5, label='
19         Area In T') # Add hatch inside the triangle
20
21 # Set axis limits
22 plt.xlim(0, 4)
23 plt.ylim(-1, 3)
24
25 # Position the legends at the top and right
26 plt.legend(loc='upper right', bbox_to_anchor=(1.2, 1.0))
27
28 # Save plot as PDF
29 plt.savefig('c.png', bbox_inches='tight')
```

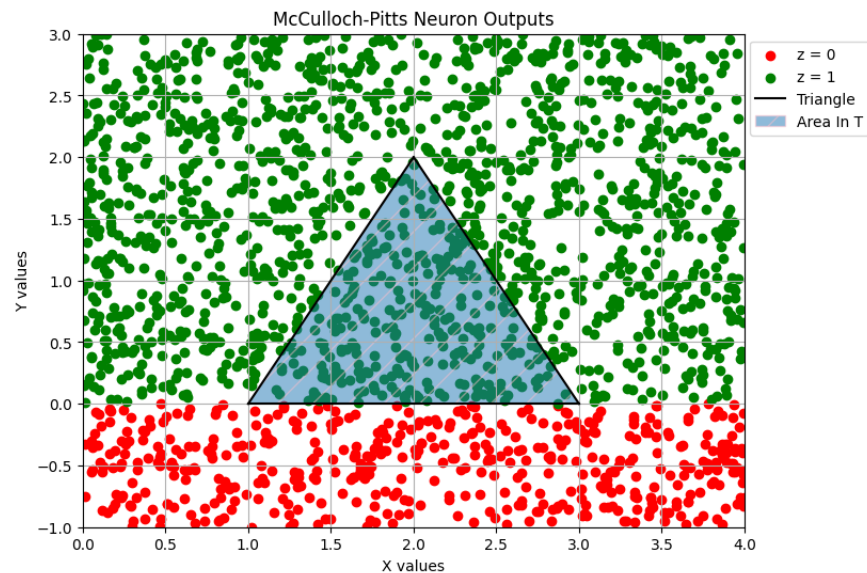


29

30 `plt.show()`

Code 7: import libraries (Python)

خروجی این طبقه‌بندی به صورت شکل **شکل ۲** آمده است. همانطور که قابل مشاهده است، چون می‌خواهیم ناحیه بالای خط جزو کلاس سبز باشد که داخل مثلث قرار بگیرد، بنابراین این کار به درستی صورت گرفته است



شکل ۲: بلوک دیاگرام مدل طبقه‌بند خطی

سپس برای خط

$$2x - y = 2$$

در کد ۱۵ در خط دوم قسمت فراخوانی، مقادیر  $[1, 2]$  را برای وزن‌دهی و مقدار ۲ را برای آستانه مقداردهی می‌کنیم. مقدار  $z_2$  را به خروجی کلاس McCulloch-Pitts برای خط دوم پس از دادن ورودی‌های مورد نظر نسبت می‌دهیم. در انتهای این تابع مقدار  $z_2$  را که خروجی تابع برای خط دوم بود را برمی‌گردانیم.

خروجی این طبقه‌بندی به صورت **شکل ۳** قابل مشاهده است.

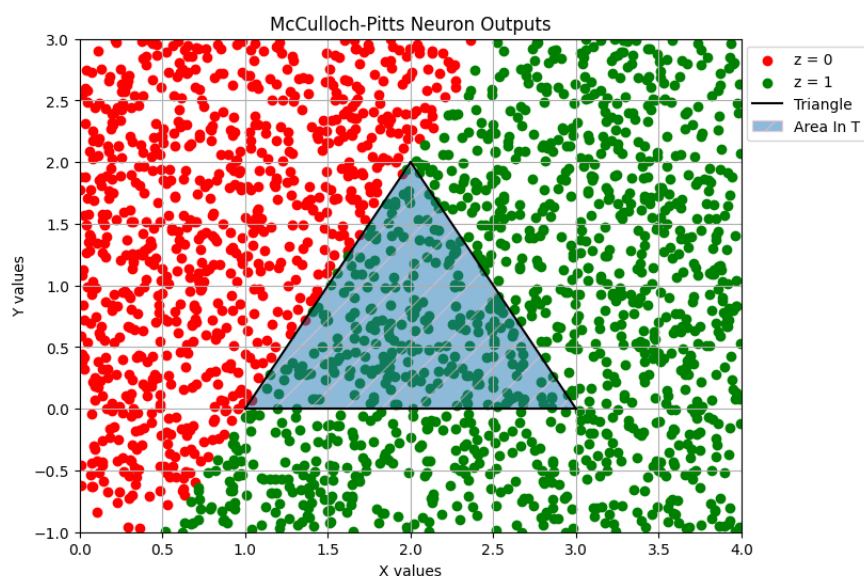
سمت چپ این خط به عنوان کلاس ۰ و سمت راست این خط به عنوان کلاس ۱ پیش‌بینی شده است که مطابق خواسته ما و ناحیه مثلث است.

در نهایت برای خط

$$-2x - y = -6$$

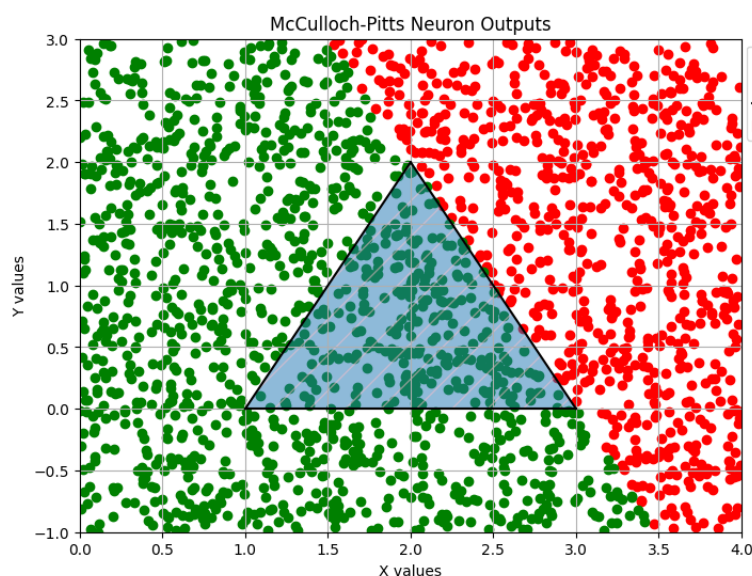
بنابراین در کد ۱۶ در خط سوم قسمت فراخوانی، مقادیر  $[1, 2]$  را برای وزن‌دهی و مقدار -۶ را برای آستانه مقداردهی می‌کنیم. مقدار  $z_3$  را به خروجی کلاس McCulloch-Pitts برای خط سوم پس از دادن ورودی‌های مورد نظر نسبت می‌دهیم. در انتهای این تابع مقدار  $z_3$  را که خروجی تابع برای خط سوم بود را برمی‌گردانیم.





شکل ۳: بلوک دیاگرام مدل طبقه‌بند خطی

خروجی طبقه‌بندی برای خط سوم به صورت زیر آمده است.



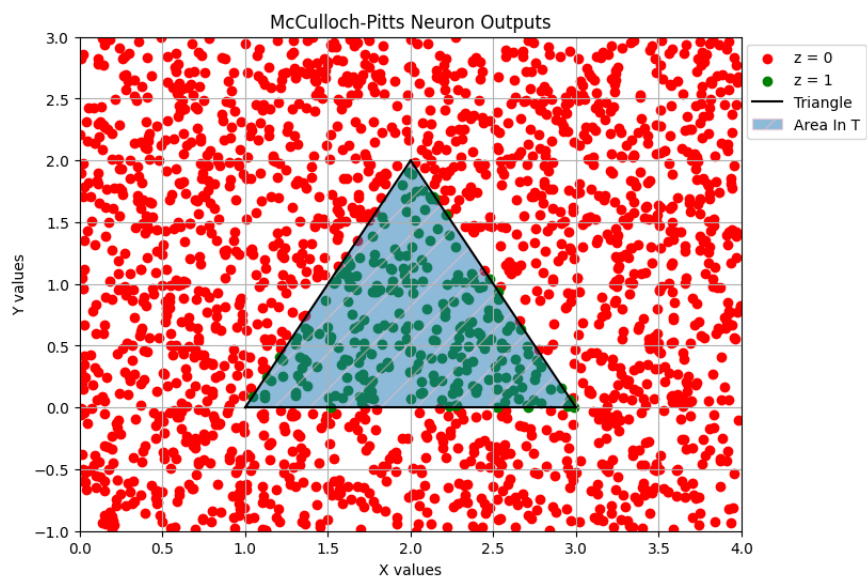
شکل ۴: بلوک دیاگرام مدل طبقه‌بند خطی

سمت راست این خط به عنوان کلاس ۰ و سمت چپ آن به عنوان کلاس ۱ طبقه‌بندی شده است که مورد نظر ماست. باید توجه داشت که اگر برعکس این اتفاق بیوفتد، می‌توان یک منفی در وزن‌ها و آستانه داده شده به نرون ضرب کرد و یا مقداردهی سبز و قرمز را به کلاس‌های طبقه‌بندی شده عوض کرد.

برای جمع‌بندی سه طبقه‌بندی انجام شده باید نتیجه هر سه طبقه‌بندی را ترکیب کرد. بنابراین در خط چهارم کد ۱۵ برای یک نرون دارای سه ورودی، وزن‌های  $[1, 1, 1]$  و آستانه ۳ را مقداردهی کرده و  $Z^4$  را به عنوان خروجی این نرون پس از دادن ورودی‌های آن (سه)



خط مورد نظر) نسبت می‌دهیم.  
خروجی این طبقه‌بندی را نیز به صورت زیر قابل مشاهده است.



شکل ۵: بلوک دیاگرام مدل طبقه‌بند خطی

در قسمت دوم این سوال خواسته شده است که اثر اضافه کردن دو تابع فعال ساز مختلف به فرآیند تصمیم‌گیری مورد بررسی قرار گیرد. تابع فعال‌ساز در مورد نورون McCulloch-Pitts که مورد آموزش قرار نمی‌گیرد و ما با توجه به استدلال و ناحیه مورد نظر مسئله مقادیر وزن و آستانه را مقداردهی می‌کنیم، زیاد مورد نیاز و جزء اساسی محسوب نمی‌شود. با این حال این کار در ادامه مورد ارزیابی قرار گرفته است.

به عنوان اولین تابع فعال‌ساز از تابع linear استفاده می‌کنیم. بیشترین میزان استفاده در این تابع فعال‌ساز در پروژه‌های مربوط به رگرسیون است اما گاهی در مسایل مربوط به طبقه‌بندی نیز مورد استفاده قرار می‌گیرند. از آنجایی که در این سوال نیز کار طبقه‌بندی صورت می‌گیرد، اما به دلیل ماهیت نورون McCulloch-Pitts و نتیجه متفاوت اضافه کردن سایر توابع فعال‌ساز، ابتدا از این تابع استفاده می‌کنیم که صرفاً خروجی نهایی هر نورون را از قسمت‌های قبلی جدا کند اما تغییری در خروجی ایجاد نکند.

به عنوان تابع فعال‌ساز دوم نیز از تابع leaky relu استفاده می‌کنیم. مزیت این تابع نسبت به ReLU این است که از آنجایی که تابع ReLU خروجی هر ورودی‌ای که مقدار منفی داشته باشد را صفر می‌کند و از طرفی این مقدار خروجی با حد آستانه مورد ارزیابی قرار می‌گیرد، بنابراین مطلوب نیست که مقادیر منفی صفر شوند. با استفاده از تابع leaky relu با ضریب آلفایی که دارد، مقدارهای منفی را صفر نمی‌کند و منفی نگه می‌دارد البته با اعمال تضعیف اندازه. این توابع به صورت زیر تعریف می‌شود.

```
1 def linear(x):
2     return x
3
4 def leaky_relu(x, alpha=0.01):
```



```
5 return max(alpha*x, x)
```

Code 8: import libraries (Python)

کدهای مربوط به تعریف کلاس نورون و تابع فراخوانی هر نورون با توجه به اضافه شدن تابع فعال‌ساز به صورت زیر تغییر پیدا کردند.

```
1 #define muculloch pitts
2 class McCulloch_Pitts_neuron_with_AF():
3     def __init__(self , weights , threshold, activation_func = None):
4         self.weights = weights      #define weights
5         self.threshold = threshold   #define threshold
6         self.activation_func = activation_func
7
8     def model(self , x):
9         #define model with threshold
10        y = self.weights @ x
11        y = y if self.activation_func == None else self.activation_func(y)
12        if y >= self.threshold:
13            return 1
14        else:
15            return 0
16
17
18 #define model for dataset
19 def Area_with_AF(x, y):
20     neur1 = McCulloch_Pitts_neuron_with_AF([0, 1], 0, linear)
21     neur2 = McCulloch_Pitts_neuron_with_AF([2, -1], 2, linear)
22     neur3 = McCulloch_Pitts_neuron_with_AF([-2, -1], -6, linear)
23     # neur3 = McCulloch_Pitts_neuron_with_AF([-2, -1], -0.06, linear)    #for
24     # leaky relu
25     neur4 = McCulloch_Pitts_neuron_with_AF([1, 1, 1], 3, linear)
26
27     z1 = neur1.model(np.array([x, y]))
28     z2 = neur2.model(np.array([x, y]))
29     z3 = neur3.model(np.array([x, y]))
30     z4 = neur4.model(np.array([z1, z2, z3]))
```



```

30
31 # return str(z1) + str(z2)
32 return list([z4])

```

Code 9: import libraries (Python)

به عنوان تابع فعال‌ساز سوم هم از sigmoid استفاده شده است. مزیت استفاده از این تابع فعال‌ساز به دلیل مپ کردن مقدار خروجی نوروں بین بازه ۰ و ۱ است. البته باید توجه داشت که باید برای هر یک از نوروں‌ها مقدار آستانه را با توجه به sigmoid تغییر داده شود. یعنی به صورت زیر مقادیر حد آستانه باید تغییر پیدا کند.

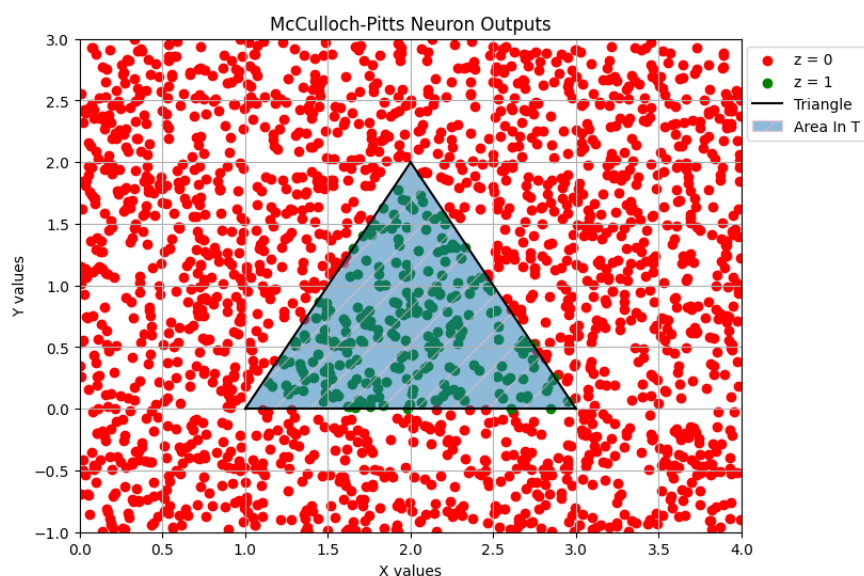
```

1 #define model for dataset
2 def Area_with_AF(x, y):
3     neur1 = McCulloch_Pitts_neuron_with_AF([0, 1], sigmoid(0), sigmoid)
4     neur2 = McCulloch_Pitts_neuron_with_AF([2, -1], sigmoid(2), sigmoid)
5     neur3 = McCulloch_Pitts_neuron_with_AF([-2, -1], sigmoid(-6), sigmoid)
6     neur4 = McCulloch_Pitts_neuron_with_AF([1, 1, 1], sigmoid(3), sigmoid)
7
8     z1 = neur1.model(np.array([x, y]))
9     z2 = neur2.model(np.array([x, y]))
10    z3 = neur3.model(np.array([x, y]))
11    z4 = neur4.model(np.array([z1, z2, z3]))
12
13    # return str(z1) + str(z2)
14    return list([z4])

```

Code 10: import libraries (Python)

خروجی شبکه با استفاده از هر یک از این توابع فعال‌ساز بصورت شکل ۶ است. مشاهده می‌شود که با استفاده از تابع فعال‌ساز linaer خروجی مورد نظر سوال که مشابه حالت قبلی نیز هست، به دست آمده است. تمامی نقاط رندوم تولید شده در داخل مثلث به رنگ سبز و متعلق به کلاس ۱ درآمده‌اند. این خروجی با همان وزن‌های استفاده شده در قسمت قبل بدست آمده است. یعنی با استفاده از این تابع فعال‌ساز که تغییری در خروجی آن نوروں ایجاد نمی‌کند، نتیجه مطلوب بدون نیاز به طراحی مجدد حاصل می‌شود. در حالت استفاده از تابع فعال‌ساز leaky relu به دلیل اینکه با توجه به آستانه در نظر گرفته شده برای خط سوم که مقداری منفی شده بود، و اینکه این تابع اندازه مقادیر منفی را تضعیف می‌کند، در حد آستانه نوروں سوم تغییری ایجاد شد که این شیفت مربوط به ضریب آلفای تابع فعال‌ساز اثر داده شود.



شکل ۶: خروجی هر سه حالت استفاده از دو تابع فعال‌ساز مختلف

## ۲ سوال دوم

۱.۲ دیتاست CWRU Bearing که در «مینی پروژه شماره یک» با آن آشنا شدید را به خاطر آورید. علاوه بر دو کلاسی که در آن مینی پروژه در نظر گرفتید، با مراجعه به صفحه داده های عیب در حالت ۱۲،  $k$  دو کلاس دیگر نیز از  $X-007@6OR$  و  $X-007B1$  اضافه کنید. با انجام این کار یک کلاس داده سالم و سه کلاس طریق فایل های  $X-007B1$  از داده های دارای سه عیب متفاوت خواهید داشت. در مورد این که هر فایل مربوط به چه نوع عیبی است به صورت کوتاه توضیح دهید. سپس در ادامه، تمام کارهایی که در بخش «۲» سوال دوم «مینی پروژه یک» برای استخراج ویژگی و آماده سازی دیتا انجام داده بودید را روی دیتاست جدید خود پیاده سازی کنید. در قسمت تقسیم بندی داده ها، یک بخش برای «اعتبارسنجی» به بخش های «آموزش» و «آزمون» اضافه کنید و توضیح دهید که کاربرد این بخش چیست.

داده ها برای یاتاقان های معمولی و عیب های در درایو تک نقطه ای و در انتهای فن جمع آوری شد. داده ها در ۱۲۰۰۰ نمونه در ثانیه و ۴۸۰۰۰ نمونه در ثانیه برای آزمایش های یاتاقان انتهایی درایو جمع آوری شد. تمام داده های بلبرینگ انتهایی فن در ۱۲۰۰۰ نمونه در ثانیه جمع آوری شد

هدف استفاده از این نوع دیتاست، تعمیر و نگهداری پیش بینی کننده در ماشین ها پیش بینی خرابی ها است. در ماشین های دوار، قطعه ای که بیشترین آسیب را متحمل می شود، بلبرینگ ها هستند. هدف اصلی این تحقیق تشخیص خرابی بلبرینگ با استفاده از حداقل مجموعه مشاهدات و انتخاب حداقل تعداد ویژگی است. مجموعه آزمایشی برای جمع آوری داده ها شامل یک موتور القایی الکتریکی ۰.۳۷ اسب بخار در یک سر، یک مبدل گشتاور در وسط، و یک دینامومتر در انتهای دیگر است که بار را شبیه سازی می کند. سنسورها، به ویژه شتاب سنج ها، روی یاتاقان های انتهای شفت موتور و روی فن داخل محفظه موتور قرار گرفتند و ۱۲۰۰۰ و ۴۸۰۰۰ نمونه در ثانیه از عیوب را جمع آوری کردند. سرعت و قدرت از طریق مبدل گشتاور بدست آمد و به صورت دستی ثبت شد.



نقص هایی با استفاده از تخلیه الکتریکی به بلبرینگ های SKF اضافه شد که باعث خرابی در تاج داخلی، توپ ها و مسیر بیرونی با قطرهای مختلف از ۰۰۷.۰ اینچ تا ۰۴.۰ اینچ شد. علاوه بر این، این آزمایش ها با تغییر سرعت چرخش و بار انجام شد. مجموعه داده شامل ۱۶۱ رکورد است که به چهار گروه تقسیم می شوند: ۴۸k خط پایه معمولی (داده های بدون خط)، ۴۸k خطای انتهای درایو، ۱۲k خطای انتهای درایو و ۱۲k خطای انتهای فن. آنها حاوی اطلاعاتی بر اساس بارها و سرعت موتور هستند. در مورد نام فایل ها، حرف اول نشان دهنده موقعیت نقص، سه عدد بعدی قطر خرابی و آخرین عدد نشان دهنده بار است. به عنوان مثال، فایل IR۰۰۷-۰ دارای اطلاعات خرابی تاج داخلی است، با قطر شکست ۰۰۷.۰ اینچ برای بار موتور ۰ اسب بخار. هر فایل داده در مجموعه داده CWRU از داده هایی با طول های مختلف تشکیل شده است و مضرب ۲ نیست، علاوه بر این مجموعه ای بزرگ، متنوع و پیچیده است. [b۴]

ابتدا مجموعه داده های مورد خواسته مسئله و کتابخانه های مورد نیاز را فراخوانی می کنیم. سپس چهار دیتاست سالم و خطا را در متغیرهای مورد نظر وارد می کنیم. با استفاده از ویژگی key می توان به سرستون های هر کدام از دیتاست ها دسترسی پیدا کرد.

```

1 %cd /content
2
3 !gdown 1VoGTdaylo4ytrvh6wMOVTxelgxWlQfvyu
4 !gdown 10IycnfVXVGsMVRQ58-Qy-14Ypv0fYY-_
5 !gdown 1u45bE0rKvKVNDkqPtmMEQLBTd1wfo6pc
6 !gdown 1EE1bGnHm5T10pgmgoI99FEQbVyBbY5Ky
7
8 from scipy.io import loadmat
9 import numpy as np
10 from scipy.stats import skew, kurtosis
11 from sklearn.utils import shuffle
12 from sklearn.model_selection import train_test_split
13
14 # Load the .mat file
15 normal_data = loadmat('/content/97.mat')
16 fault_data_IR007 = loadmat('/content/105.mat')
17 fault_data_BB007 = loadmat('/content/118.mat')
18 fault_data_OR007 = loadmat('/content/130.mat')
19
20 print(normal_data.keys())
21 print(fault_data_IR007.keys())
22 print(fault_data_BB007.keys())
23 print(fault_data_OR007.keys())

```

Code 11: import dataset and library

خروجی سرستون های دیتاست ها بصورت شکل ۷ است.



```
dict_keys(['_header_', '_version_', '_globals_', 'X097_DE_time', 'X097_FE_time', 'X097RPM'])
dict_keys(['_header_', '_version_', '_globals_', 'X105_DE_time', 'X105_FE_time', 'X105_BA_time', 'X105RPM'])
dict_keys(['_header_', '_version_', '_globals_', 'X118_DE_time', 'X118_FE_time', 'X118_BA_time', 'X118RPM'])
dict_keys(['_header_', '_version_', '_globals_', 'X130_DE_time', 'X130_FE_time', 'X130_BA_time', 'X130RPM'])
```

شکل ۷: نتیجه آموزش و ارزیابی با استفاده از SGD

سپس از مجموعه داده سالم ستون X097-DE-time و از مجموعه داده های خطا ستون های مربوط به time را به عنوان نمونه با توجه به خواسته مسئله انتخاب می کنیم و در متغیرهای مناسب می ریزیم. برای داشتن دید مناسب از ابعاد مجموعه داده ابتدا نوع دیتای ستون های مورد نظر و همچنین ابعاد آن ها را خروجی می گیریم.

```
1 normal_data_variable = normal_data['X097_DE_time']
2 print(type(normal_data_variable))
3 print(normal_data_variable.shape)
4
5 fault_data_IR007_variable = fault_data_IR007['X105_DE_time']
6 print(fault_data_IR007_variable.shape)
7
8 fault_data_BB007_variable = fault_data_BB007['X118_DE_time']
9 print(fault_data_BB007_variable.shape)
10
11 fault_data_OR007_variable = fault_data_OR007['X130_DE_time']
12 print(fault_data_OR007_variable.shape)
```

Code 12: import dataset and library

خروجی مورد نظر به صورت شکل ۸ است.

```
<class 'numpy.ndarray'>
(243938, 1)
(121265, 1)
(122571, 1)
(121991, 1)
```

شکل ۸: نتیجه آموزش و ارزیابی با استفاده از SGD

در این با انتخاب تعداد نمونه برابر ۱۲۰ و اندازه طول ۲۲۰ از هر کلاس جداسازی شده است. این ۱۲۰ نمونه با طول ۲۲۰ به صورت پشت سر هم قرار گرفته اند. در نهایت ابعاد هر کدام برابر (۱، ۲۶۴۰۰) می شود. سپس برا تکمیل ابعاد مجموعه داده با استفاده از دستور reshape، دیتاها به صورتی درمی آیند که سطرهای آن دربرگیرنده ۱۲۰ نمونه مورد خواسته مسئله باشد.





```
1 # Extract 10 samples, each containing 5 data points
2 n_samples = 120
3 n_data = 220
4 ext_normal_data = normal_data_variable[:n_samples * n_data, 0]
5 ext_fault_data_IR007 = fault_data_IR007_variable[:n_samples * n_data, 0]
6 ext_fault_data_BB007 = fault_data_BB007_variable[:n_samples * n_data, 0]
7 ext_fault_data_OR007 = fault_data_OR007_variable[:n_samples * n_data, 0]
```

Code 13: import dataset and library

```
1 # Reshape the extracted data to have 10 rows and 5 columns
2 ext_normal_data = ext_normal_data.reshape(n_samples, n_data)
3 ext_fault_data_IR007 = ext_fault_data_IR007.reshape(n_samples, n_data)
4 ext_fault_data_BB007 = ext_fault_data_BB007.reshape(n_samples, n_data)
5 ext_fault_data_OR007 = ext_fault_data_OR007.reshape(n_samples, n_data)
6
7
8 print(ext_normal_data.shape)
9 print(ext_fault_data_IR007.shape)
10 print(ext_fault_data_BB007.shape)
11 print(ext_fault_data_OR007.shape)
```

Code 14: import dataset and library

خروجی این کد نیز بصورت زیر است.

```
(120, 220)
(120, 220)
(120, 220)
(120, 220)
```

شکل ۹: نتیجه آموزش و ارزیابی با استفاده از SGD

با توجه به جدول مسئله تعداد هشت عدد از روش‌های ذکر شده انتخاب شد. با استفاده از این روش‌ها ویژگی‌های دیتاست مورد استخراج می‌شود.

تابع feature-extraction طوری کدنویسی شده است که بتواند محاسبه هر کدام از این روش‌ها و اعمال بر روی دیتاست‌های ورودی را انجام دهد. برای محاسبه هر کدام از این روش‌ها از کتابخانه numpy استفاده شده است. قسمت `axis=1` باعث می‌شود که الگوریتم





مورد نظر بر روی سطرها که نشان‌دهنده نمونه‌های دیتاست هستند اعمال شود. در قسمت بعدی خروجی ویژگی‌های استخراج‌شده با دستور column-stack بصورت ستونی به یکدیگر متصل شدند تا هر دیتاست با ابعاد (۸، ۱۲۰) ساخته شود. دی قسمت پایانی نیز یک ستون به دیتاست اضافه شده است که نشان‌دهنده برچسب کلاس است. این ستون با توجه به نوع کلاس از ۰ تا ۳ متغیر است.

```
1 def feature_extraction (data, class_num):
2     #Normal data feature extraction
3     standard_deviations = np.std(data, axis=1)
4     skewnesses = skew(data, axis=1)
5     kurtoses = kurtosis(data, axis=1)
6     peak_to_peaks = np.ptp(data, axis=1)
7     root_mean_squares = np.sqrt(np.mean(np.square(data), axis=1))
8     means = np.mean(data, axis=1)
9     absolute_means = np.mean(np.abs(data), axis=1)
10    peaks = np.max(data, axis=1)
11
12    print("Standard Deviations:", standard_deviations.shape)
13    print("Skewnesses:", skewnesses.shape)
14    print("Kurtoses:", kurtoses.shape)
15    print("Peak to Peaks:", peak_to_peaks.shape)
16    print("Root Mean Squares:", root_mean_squares.shape)
17    print("Means:", means.shape)
18    print("Absolute Means:", absolute_means.shape)
19    print("Peaks:", peaks.shape)
20
21    ext_feature_dataset = np.column_stack((standard_deviations, skewnesses,
22                                           kurtoses, peak_to_peaks,
23                                           root_mean_squares, means, absolute_means,
24                                           peaks))
25
26    ext_feature_dataset = np.hstack((ext_feature_dataset, class_num * np.ones((
27        len(ext_feature_dataset), 1))))
28    print("final dataset :", ext_feature_dataset.shape)
```



```
29 return ext_feature_dataset
```

Code 15: import dataset and library

در ادامه برای استفاده از این تابع برای هر کدام از دیتاهای نرمال و خطا، کد زیر به همراه خروجی (مشابه برای هر چهار نوع دیتا) قایل مشاهده است.

```
1 normal_ext_feature_dataset = feature_extraction(ext_normal_data, 0)
2 IR007_ext_feature_dataset = feature_extraction(ext_fault_data_IR007, 1)
3 BB007_ext_feature_dataset = feature_extraction(ext_fault_data_BB007, 2)
4 OR007_ext_feature_dataset = feature_extraction(ext_fault_data_OR007, 3)
5
6
7 Standard Deviations: (120,)
8 Skewnesses: (120,)
9 Kurtoses: (120,)
10 Peak to Peaks: (120,)
11 Root Mean Squares: (120,)
12 Means: (120,)
13 Absolute Means: (120,)
14 Peaks: (120,)
15 final_dataset : (120, 9)
```

Code 16: import dataset and library

در نهایت برای این که دیتاست تشکیل شده مانند یک دیتاستی که دارای کلاس‌های مختلف است به صورت استاندارد دربیاید، چهار ماتریس تشکیل شده با هم ادغام شده است. کد و خروجی آن بصورت زیر قابل بررسی است.

```
1 final_data = np.concatenate((normal_ext_feature_dataset ,
    IR007_ext_feature_dataset , BB007_ext_feature_dataset ,
    OR007_ext_feature_dataset), axis=0)
2 final_data.shape
```

Code 17: import dataset and library

در این قسمت پس از مخلوط کردن دیتاها با دستور shuffle، دیتاهای مربوط به ویژگی‌های دیتاست (همه ستون‌ها به جز ستون آخر) به متغیر X و دیتای برچسب به متغیر y تعلق پیدا می‌کند. در نهایت دیتاها به دو دسته آموزش و ارزیابی تقسیم می‌شوند. نسبت داده‌های



(480, 9)

شکل ۱۰: نتیجه آموزش و ارزیابی با استفاده از SGD

آموزش به ارزیابی در این قسمت با توجه به اندازه کلی دیتاست برابر ۱۰ درصد در نظر گرفته شده است. در سوال خواسته شده است که یک بخش اعتبارسنجی نیز در نظر گرفته شود. این بخش در ادامه از قسمت دیتاست آموزش جدا خواهد شد. تقسیم به بخش اعتبارسنجی به عنوان یک گام اساسی در یادگیری ماشین شناخته می شود. معمول است که این بخش باید از دیتاست آموزش جدا شود، یعنی پس از اینکه دیتاست تست را از آموزش جدا کردیم، از دیتاست باقی مانده، بخش اعتبارسنجی را جدا می کنیم. عدم اشتراک داده در این سه بخش بسیار مهم است. از دلایل اهمیت و کاربردهای این بخش می توان به این موارد اشاره کرد: ارزیابی مدل: از داده های اعتبارسنجی برای ارزیابی مدل هنگام آموزش استفاده می شود. کمک به تیون کردن هایپرپارامترها و تشخیص اینکه کدام مدل با چه پارامترهایی مناسب تر است، از مزایای استفاده از این بخش است. جلوگیری از اورفیت و آندرفیت شدن مدل: با ارزیابی نمودارهای مربوط به loss داده های اعتبارسنجی و مقایسه با نمودار آموزش، می توان پیش از تمام شدن کل آموزش به ایجاد شدن این مشکلات پی برد و مشکلات را حل کرد. توقف زودهنگام: با مقایسه نمودارهای آموزش و اعتبارسنجی می توان زمانی را که مقدار loss آموزش و اعتبارسنجی از حدی کمتر نمی شوند را تشخیص داده و از آموزش بیشتر و صرف هزینه زمانی و سخت افزاری اضافه جلوگیری کرد. کد مورد نظر و خروجی نهایی بصورت زیر است.

```
1 shuffled_final_data = shuffle(final_data)
2
3 X = shuffled_final_data[:, :-1]
4 y = shuffled_final_data[:, -1]
5 print("X :", X.shape)
6 print("y :", y.shape)
7
8 x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.1)
9 x_train.shape, x_test.shape, y_train.shape, y_test.shape
```

Code 18: import dataset and library

((432, 8), (48, 8), (432,), (48,))

شکل ۱۱: نتیجه آموزش و ارزیابی با استفاده از SGD

اما این ابعاد برای برچسب ها در ادامه مشکل ساز خواهد بود بنابراین نیاز است که با تغییر ابعاد، آن را به صورتی در بیاوریم که بعد ۱ در قسمت ستون وارد شود. با استفاده از دستور reshape این کار را انجام می دهیم.

```
1 # Reshape y_train and y_test
```



```

2 y_train = np.reshape(y_train, (-1, 1))
3 y_test = np.reshape(y_test, (-1, 1))
4 x_train.shape, x_test.shape, y_train.shape, y_test.shape

```

Code 19: import dataset and library

در نهایت ابعاد دیتاهای مورد نظر به صورت زیر در می آید.

```
((432, 8), (48, 8), (432, 1), (48, 1))
```

شکل ۱۲: نتیجه آموزش و ارزیابی با استفاده از SGD

نیاز است که فرایند نرمال سازی پس از تقسیم داده ها به آموزش و تست صورت بگیرد چرا که داده های تست نباید توسط مدل دیده شوند و با اعمال ویژگی آن ها در فرایند نرمال سازی، برخی ویژگی ها وارد مدل می شود و ارزیابی صورت گرفته درست نخواهد بود. در نتیجه فرایند نرمال سازی با استفاده از اطلاعات داده ها آموزش بر روی داده های آموزش و ارزیابی صورت می پذیرد. برای نرمال سازی داده ها در این سوال از روش MinMaxScaler استفاده می کنیم. کد مورد استفاده بصورت زیر است.

```

1 from sklearn.preprocessing import MinMaxScaler, StandardScaler
2
3 # Initialize MinMaxScaler
4 scaler = MinMaxScaler()
5
6 scaler.fit(x_train)
7 x_train = scaler.transform(x_train)
8 x_test = scaler.transform(x_test)

```

Code 20: import dataset and library

۲.۲ یک مدل MLP ساده با ۲ لایه پنهان یا بیش تر بسازید. بخشی از داده های آموزش را برای اعتبارسنجی کنار بگذارید و با انتخاب بهینه ساز و تابع اتلاف مناسب، مدل را آموزش دهید. نمودارهای اتلاف و Accuracy مربوط به آموزش و اعتبارسنجی را رسم و نتیجه را تحلیل کنید. نتیجه تست مدل روی داد های آزمون را با استفاده ماتریس درهم ریختگی و report-classification نشان داده و نتایج به صورت دقیق تحلیل کنید.

ابتدا لازم است کتابخانه های مورد نیاز را فراخوانی کنیم.



```
1 from keras.utils import to_categorical
2 from sklearn.preprocessing import LabelBinarizer
3
4 from sklearn.neural_network import MLPClassifier, MLPRegressor
5 from imblearn.under_sampling import RandomUnderSampler
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8 import numpy as np
9
10 import tensorflow as tf
11 from tensorflow import keras
12 from keras import preprocessing
13 from keras.models import Sequential
14 from keras.layers import Dense
```

Code 21: import dataset and library

کار طبقه‌بندی در این سوال را با استفاده از یک شبکه MLP با دو لایه پنهان انجام می‌دهیم. از فریم‌ورک کراس استفاده شده است. جزئیات شبکه دوباره به صورت زیر قابل مشاهده است. با دستور Sequential مدلی فراخوانی می‌شود که می‌توان در آن لایه‌های مورد نظر شبکه را در آن وارد کرده و تعریف کرد. سپس در دو نوبت با استفاده از دستور add لایه پنهان اولی با ۱۵ نورون و لایه دوم را با ۸ نورون ایجاد می‌کنیم. با توجه به اینکه تعداد کلاس‌های دیتاست ۴ تاست، باید لایه طبقه‌بندی آخر را با ۴ نورون ایجاد کنیم. با استفاده از دستور summery می‌توان جزئیات شبکه را مشاهده کرد. تابع فعال‌ساز لایه‌های پنهان ReLU و لایه طبقه‌بند softmax در نظر گرفته شده‌اند.

```
1 model_2 = Sequential()
2
3 # Add the first hidden layer with 50 neurons and linear activation function
4 model_2.add(Dense(50, activation='relu', input_shape=(x_train.shape[1],)))
5
6 # Add the second hidden layer with 30 neurons and linear activation function
7 model_2.add(Dense(30, activation='relu'))
8
9 # Add an output layer with 1 neuron and linear activation function
10 model_2.add(Dense(4, activation='softmax'))
11
12 model_2.summary()
```

Code 22: import dataset and library



```
Model: "sequential_1"
Layer (type)                Output Shape                Param #
=====
dense_3 (Dense)              (None, 15)                  135
dense_4 (Dense)              (None, 8)                   128
dense_5 (Dense)              (None, 4)                   36
=====
Total params: 299 (1.17 KB)
Trainable params: 299 (1.17 KB)
Non-trainable params: 0 (0.00 Byte)
```

شکل ۱۳: نتیجه آموزش و ارزیابی با استفاده از SGD

با استفاده از کد ۳ می‌توان شبکه MLP را آموزش داد. با دستور `compile` و با در نظر گرفتن بهینه‌ساز Adam و تابع اتلاف Spar-`seCategoricalCrossentropy` به تعداد ۱۰۰ اپیاک مدل را آموزش می‌دهیم. مقدار بخش اعتبارسنجی را برابر ۲۰ درصد قرار می‌دهیم. مقدار بچ‌سایز را هم برابر ۱۰ قرار می‌دهیم. با این کار داده‌ها در بسته‌های ۱۰ تایی مورد آموزش قرار می‌گیرند.

```
1 model_2.compile(optimizer='adam', loss='SparseCategoricalCrossentropy',
    metrics=['accuracy'])
2 history = model_2.fit(x_train, y_train, validation_split=0.2, epochs=100 ,
    batch_size=10)
```

Code 23: import dataset and library

پس از پایان آموزش باید فرآیند آموزش را مورد ارزیابی قرار داد. یک روش معمول و مناسب برای اینکار رسم و مقایسه نمودارهای اتلاف و accuracy داده‌های آموزش و اعتبارسنجی است.

```
1 import matplotlib.pyplot as plt
2
3 plt.plot(history.history['loss'], label='train')    # Training loss
4 plt.plot(history.history['val_loss'], label='val')  # Validation loss
5
6 plt.legend(['Training Loss', 'Validation Loss'])
7 plt.xlabel("Epochs")
8 plt.ylabel("Loss")
9 plt.show()
```



```
Epoch 1/100
35/35 [=====] - 1s 9ms/step - loss: 1.4030 - accuracy: 0.2522 - val_loss: 1.3863 - val_accuracy: 0.2529
Epoch 2/100
35/35 [=====] - 0s 3ms/step - loss: 1.3588 - accuracy: 0.2522 - val_loss: 1.3449 - val_accuracy: 0.2989
Epoch 3/100
35/35 [=====] - 0s 3ms/step - loss: 1.3188 - accuracy: 0.4609 - val_loss: 1.2926 - val_accuracy: 0.4943
Epoch 4/100
35/35 [=====] - 0s 3ms/step - loss: 1.2701 - accuracy: 0.5304 - val_loss: 1.2431 - val_accuracy: 0.5287
Epoch 5/100
35/35 [=====] - 0s 3ms/step - loss: 1.2058 - accuracy: 0.5362 - val_loss: 1.1662 - val_accuracy: 0.5747
Epoch 6/100
35/35 [=====] - 0s 3ms/step - loss: 1.1059 - accuracy: 0.5391 - val_loss: 1.0537 - val_accuracy: 0.5517
Epoch 7/100
35/35 [=====] - 0s 4ms/step - loss: 0.9851 - accuracy: 0.5710 - val_loss: 0.9337 - val_accuracy: 0.5632
Epoch 8/100
35/35 [=====] - 0s 3ms/step - loss: 0.8679 - accuracy: 0.5855 - val_loss: 0.8217 - val_accuracy: 0.6552
Epoch 9/100
35/35 [=====] - 0s 4ms/step - loss: 0.7654 - accuracy: 0.6522 - val_loss: 0.7251 - val_accuracy: 0.7356
Epoch 10/100
35/35 [=====] - 0s 3ms/step - loss: 0.6766 - accuracy: 0.7333 - val_loss: 0.6381 - val_accuracy: 0.8276
Epoch 11/100
35/35 [=====] - 0s 3ms/step - loss: 0.5991 - accuracy: 0.8145 - val_loss: 0.5618 - val_accuracy: 0.8851
Epoch 12/100
35/35 [=====] - 0s 4ms/step - loss: 0.5305 - accuracy: 0.9246 - val_loss: 0.4933 - val_accuracy: 0.9080
Epoch 13/100
35/35 [=====] - 0s 3ms/step - loss: 0.4709 - accuracy: 0.9507 - val_loss: 0.4385 - val_accuracy: 0.9425
Epoch 14/100
35/35 [=====] - 0s 3ms/step - loss: 0.4210 - accuracy: 0.9768 - val_loss: 0.3897 - val_accuracy: 0.9655
Epoch 15/100
35/35 [=====] - 0s 4ms/step - loss: 0.3760 - accuracy: 0.9797 - val_loss: 0.3511 - val_accuracy: 0.9885
Epoch 16/100
```

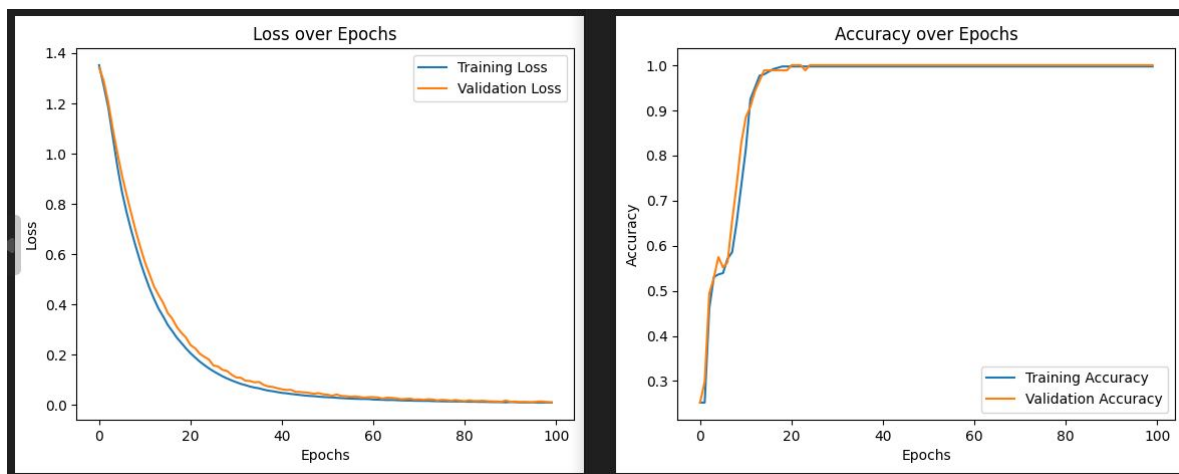
شکل ۱۴: نتیجه آموزش و ارزیابی با استفاده از SGD

```
10
11 plt.plot(history.history['accuracy'], label='Training Accuracy')
12 plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
13 plt.xlabel('Epochs')
14 plt.ylabel('Accuracy')
15 plt.title('Accuracy over Epochs')
16 plt.legend()
17
18 plt.show()
```

Code 24: import dataset and library

نمودارهای مربوط به اتلاف و accuracy بصورت **شکل ۱۵** است. همانطور که قابل مشاهده است نمودار اتلاف در دوره‌های آموزش مختلف فایل هرگونه علامتی از اورفیت یا اندرفیت است. نمودار اعتبارسنجی فرآیند کاهشی خود را همگام با نمودار آموزش طی کرده است و تا ایپاک حدود ۴۰ روند کاهشی با شیب مناسبی داشته و پس از آن میزان کاهش بسیار کم شده است. برای نمودار Accuracy نیز تا ایپاک ۱۸ نمودار افزایشی است و سپس به میزان ۱ میرسد که این نیز نشانه آموزش مناسب است. می‌توان نتیجه گرفت که نیازی به آموزش تا ۱۰۰ دوره آموزشی نیست و با توجه به عدم بهبود نمودارها و رسیدن به مقدار نهایی، می‌توان توقف زود هنگام را اجرایی کرد. در گام بعدی باید نتیجه تست روی داده‌های ارزیابی مورد بررسی قرار گیرد. جهت انجام این کار از کد زیر استفاده می‌کنیم. ابتدا داده‌های مربوط به تست را به مدل جهت انجام طبقه‌بندی می‌دهیم. سپس با توجه به احتمالاتی بودن پیش‌بینی مدل، آن طبقه‌ای که ضریب احتمال بالاتری دارد را به عنوان پیش‌بینی نهایی در نظر می‌گیریم. نتیجه تست روی داده‌های ارزیابی به صورت **شکل ۱۶** است.

```
1 y_pred_2 = model_2.predict(x_test)
```



شکل ۱۵: نتیجه آموزش و ارزیابی با استفاده از SGD

```
2
3 from sklearn.metrics import accuracy_score
4
5 # Assuming y_test contains the true labels
6 y_pred_classes = np.argmax(y_pred_2, axis=1) # Convert predicted
   probabilities to class labels
7 accuracy = accuracy_score(y_test, y_pred_classes)
8 print("Accuracy:", accuracy)
```

Code 25: import dataset and library

```
2/2 [=====] - 0s 4ms/step
Accuracy: 1.0
```

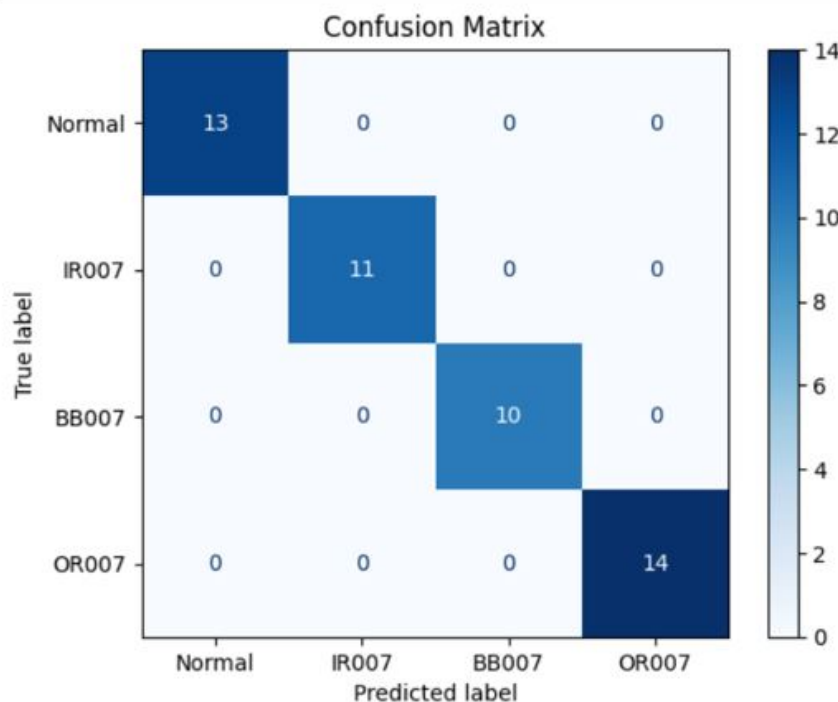
شکل ۱۶: نتیجه آموزش و ارزیابی با استفاده از SGD

با استفاده از classification-report و ماتریس درهم‌ریختگی برای داده‌های ارزیابی به بررسی نتایج دقیق‌تر می‌پردازیم. تعداد ۴۸ داده از ۴ کلاس مورد ارزیابی قرار گرفت که در تمامی کلاس‌ها پیش‌بینی به درستی اتفاق افتاده است. ماتریس درهم‌ریختگی نیز تمامی داده‌ها روی قطر اصلی هستند و پیش‌بینی اشتباهی صورت نگرفته است.





	precision	recall	f1-score	support
0.0	1.00	1.00	1.00	13
1.0	1.00	1.00	1.00	11
2.0	1.00	1.00	1.00	10
3.0	1.00	1.00	1.00	14
accuracy			1.00	48
macro avg	1.00	1.00	1.00	48
weighted avg	1.00	1.00	1.00	48



شکل ۱۷: نتیجه آموزش و ارزیابی با استفاده از SGD

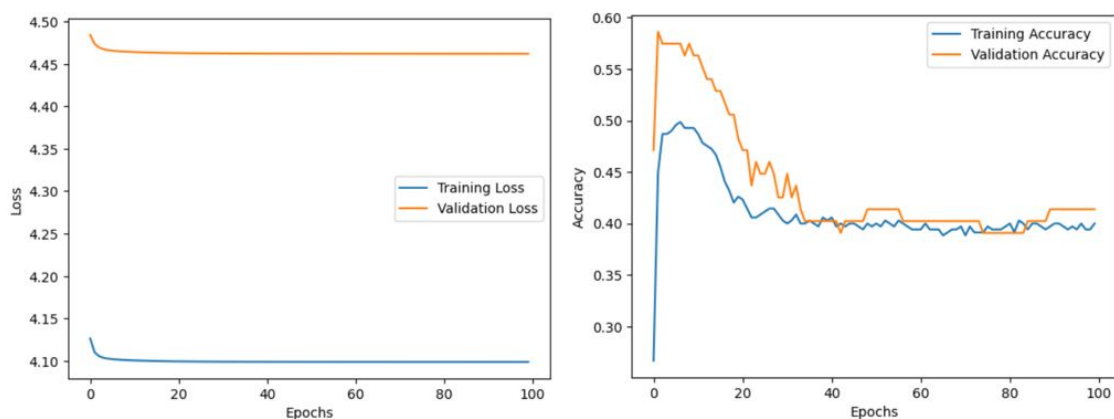
۳.۲ فرآیند سوال قبل را با یک بهینه ساز و تابع اتلاف جدید انجام داده و نتایج را مقایسه و تحلیل کنید. بررسی کنید که آیا تغییر تابع اتلاف می تواند در نتیجه اثرگذار باشد؟

برای اینکه بتوان مقایسه مناسبی بین بهینه سازها و توابع اتلاف داشت، تمامی پارامترهای قسمت قبل را برای آموزش روی بهینه ساز و تابع اتلاف جدید بدون تغییر باقی می گذاریم. در این قسمت از بهینه ساز SGD و تابع اتلاف KLDivergence استفاده می کنیم. باید توجه داشت که تابع اتلاف SparseCategoricalCrossentropy که در قسمت قبل استفاده شد، معمول تر و مورد استفاده تر است و تابع اتلاف جدید ممکن است به مطلوبی قبلی نباشد.

با آموزش مجدد به نمودارهای اتلاف و معیار شکل ۱۸ می رسمیم. همانطور که قابل مشاهده است، نمودارهای اتلاف در مرحله آموزش نتوانسته است به خوبی کاهش پیدا کند و به طور کلی آموزش مناسبی صورت نگرفته است. نمودار اعتبارسنجی نیز با فاصله زیاد و کاهش کم مواجه شده است. نمودار معیار نیز پس از افزایش نسبی در همان ایپاک های اولیه افتاده است. همه چی از آموزش نامناسب خبر می دهد



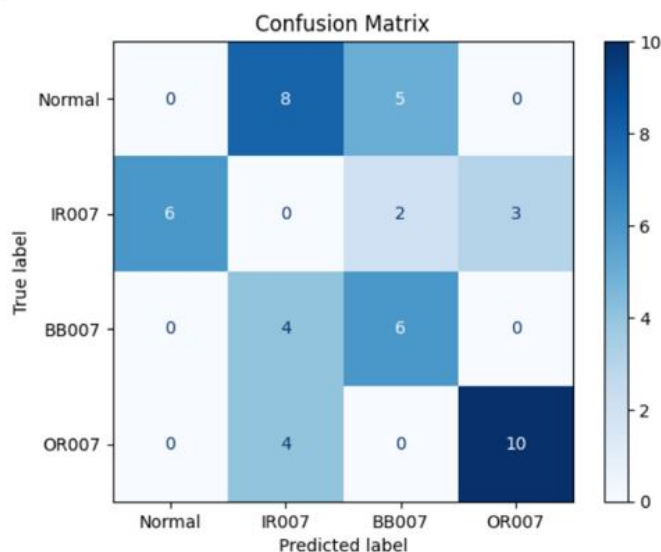
که یا ناشی از بهینه‌ساز است یا تابع اتلاف.



شکل ۱۸: نتیجه آموزش و ارزیابی با استفاده از SGD

در این حالت معیارها در شکل ۱۹ در کلاس‌های مختلف نیز نشان‌دهنده آموزش نامناسب است.

	precision	recall	f1-score	support
0.0	0.00	0.00	0.00	13
1.0	0.00	0.00	0.00	11
2.0	0.46	0.60	0.52	10
3.0	0.77	0.71	0.74	14
accuracy			0.33	48
macro avg	0.31	0.33	0.32	48
weighted avg	0.32	0.33	0.32	48

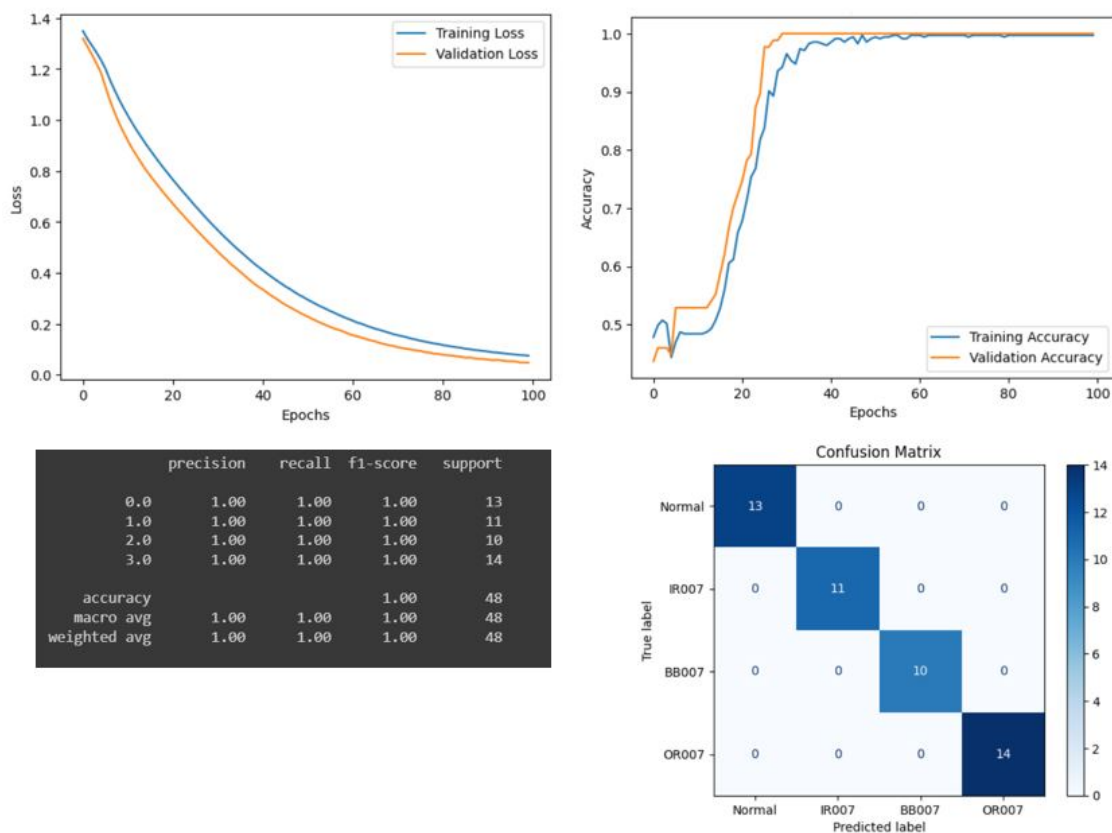


شکل ۱۹: نتیجه آموزش و ارزیابی با استفاده از SGD

در ادامه برای اینکه بفهمیم مشکل از کجا بوده است، تابع اتلاف را به همان SparseCategoricalCrossentropy برمی‌گردانیم اما بهینه‌ساز را مجدداً SGD انتخاب می‌کنیم. در این حالت نمودارها و معیارها به صورت شکل ۲۰ درآمده‌اند. نمودار اتلاف به خوبی



کاهش پیدا کرده و نمودار معیارهم تا مقدار ۱ افزایشی است. ماتریس درهم‌ریختگی و معیارهای دیگر نیز همگی نشان از آموزش مناسب می‌دهند.



شکل ۲۰: نتیجه آموزش و ارزیابی با استفاده از SGD

نتیجه‌گیری نهایی برای این قسمت را می‌توان اینگونه مطرح کرد که تغییر هم بهینه‌ساز و هم تابع اتلاف می‌تواند اثر متفاوتی در آموزش و عملکرد مدل بگذارند اما انتخاب تابع اتلاف اشتباه می‌تواند حتی در یک مجموعه داده بسیار ساده مانند این مثال، بسیار باعث افت معیارهای عملکردی شود. در آموزش سوم نیز مشاهده می‌شود که بهینه‌ساز adam سرعت همگرایی بیشتر نسبت به SGD داشته است و چون مجموعه داده انتخابی آسان بوده است، هر دو به مقدار ۱۰۰ درصد رسیده‌اند.

۴.۲ در مورد Cross-validation K-Fold Stratified و Cross-validation K-Fold و مزایای هر یک توضیح دهید. سپس با ذکر دلیل، یکی از این روش‌ها را انتخاب کرده و بخش «۲» سوال سوم را با آن پیاده سازی کنید و نتایج خود را تحلیل کنید.

روش Cross-validation K-Fold تکنیکی است که برای ارزیابی عملکرد و تعمیم‌پذیری یک مدل یادگیری ماشین استفاده می‌شود. مجموعه داده به K زیر مجموعه با اندازه مساوی تقسیم می‌شود. ابتدا دیتاست بصورت تصادفی به K بخش با اندازه مساوی تقسیم می‌شود برای هر یک از این K بخش، یک تقسیم‌بندی K-1 برای داده‌های آموزش صورت می‌گیرد. به این صورت که آن بخش به عنوان داده اعتبارسنجی استفاده می‌شود و K-1 بخش دیگر به عنوان داده آموزش. این فرآیند برای همه K بخش اجرایی می‌شود. معیار نهایی



مانند accuracy بصورت میانگینی از K تکرار محاسبه می‌شود تا یک برآورد کلی از عملکرد مدل و دیتاست به دست آید. این روش اطمینان حاصل می‌کند که یک دیتا حتما یکبار هم در دیتای آموزش باشد و هم در دیتای اعتبارسنجی. که به طور کلی منجر به پیش‌بینی قابل اعتمادتری از عملکرد مدل در مقایسه با تقسیم‌بندی معمولی و یکبار می‌شود.

روش Stratified K-Fold Cross-validation نوعی از K-Fold Cross-validation است که تضمین می‌کند هر بخش نماینده توزیع کلی کلاس‌ها است. این ویژگی در دیتاست‌هایی که عدم تعادل کلاسی دارند و یک یا چند کلاس کمتر تکرار شده‌اند مناسب است. فرآیند مانند روش قبلی است با این تفاوت که هنگام تقسیم داده به K بخش، اطمینان حاصل می‌شود که توزیع کلاس در هر بخش منعکس‌کننده توزیع کلاس در کل مجموعه داده باشد. مزیت استفاده از این روش این است که نسبت کلاس‌ها را در هر قسمت حفظ می‌کند، که منجر به تخمین‌های عملکرد قابل اعتمادتر و سازگارتر، به‌ویژه برای مجموعه داده‌های نامتعادل می‌شود. بنابراین با توجه به اینکه تعادل کلاس‌ها در مجموعه داده به چه صورت است می‌توان از یک از این دو روش استفاده کرد.

در این قسمت ابتدا نیاز است با استفاده از کتابخانه sklearn، دستور استفاده از k-fold را فراخوانی کنیم. سپس مدلی که در قسمت‌های قبلی مورد استفاده قرار داده بودیم را به شکل تابع درآوریم که در طول آموزش مدل طی دفعات مختلف بشود آن را فراخوانی کرد. مقدار k را هم با توجه به تعداد داده‌ها برابر ۵ قرار داده‌ایم.

```
1 from sklearn.model_selection import KFold
2
3 kf = KFold(n_splits=5, shuffle=True, random_state=4)
4 accuracies = []
5
6
7 def build_model():
8     model_2_3 = Sequential()
9     model_2_3.add(Dense(15, activation='relu', input_shape=(x_train.shape[1],)
10     ))
11     model_2_3.add(Dense(8, activation='relu'))
12     model_2_3.add(Dense(4, activation='softmax'))
13     model_2_3.compile(optimizer='adam', loss='SparseCategoricalCrossentropy',
14     metrics=['accuracy'])
15     return model_2_3
```

Code 26: import dataset and library

سپس با استفاده از مجموعه داده آموزش که قبلا از داده ارزیابی جدا شده بود، با استفاده از دستور kf.split و طی ۵ حلقه، آموزش مدل روی دیتاست را انجام می‌دهیم. باید توجه داشت که اینجا مجموعه داده آموزش اولیه به دو قسمت آموزش و اعتبارسنجی جدید تقسیم شده است و برای امکان ارزیابی دقیق‌تر با قسمت‌های قبلی قسمت دیتاست ارزیابی مشابه با مراحل قبلی تمرین است. کد مورد استفاده به صورت ۳۴ است.

```
1 for train_index, valid_index in kf.split(x_train):
```

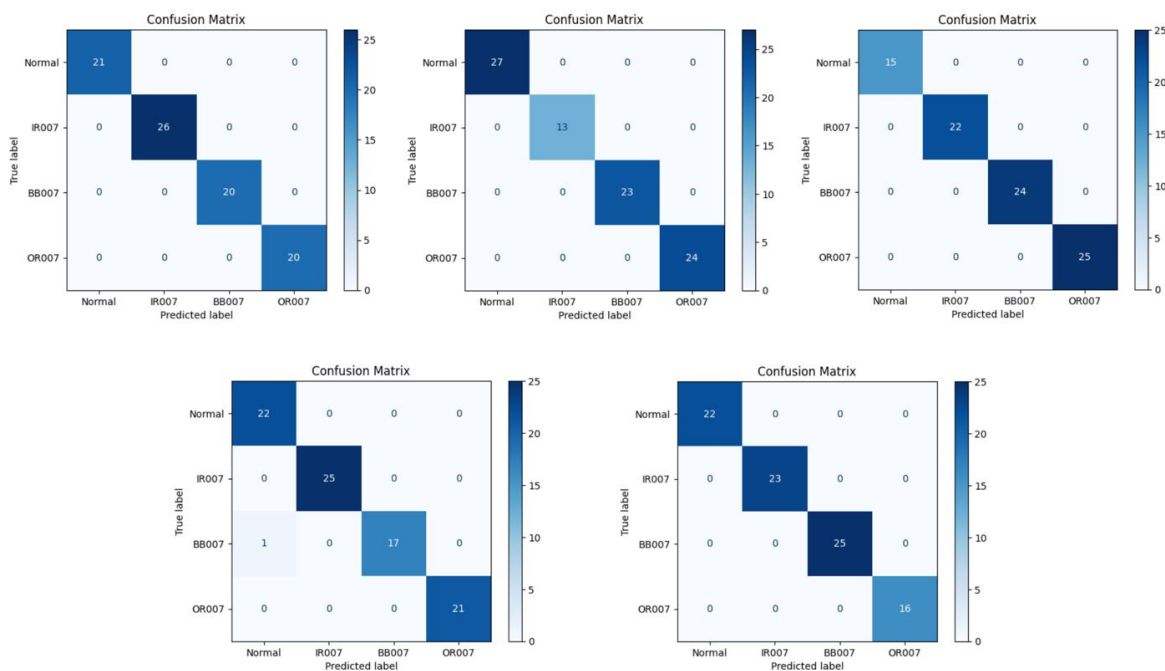


```
2 # Split the data
3 x_train_k, x_valid_k = x_train[train_index], x_train[valid_index]
4 y_train_k, y_valid_k = y_train[train_index], y_train[valid_index]
5
6 # Initialize and train the model (RandomForestClassifier in this case)
7 # model = RandomForestClassifier(n_estimators=100, random_state=42)
8 model_2_3 = build_model()
9 model_2_3.fit(x_train_k, y_train_k, validation_data=(x_valid_k, y_valid_k)
, epochs=100, batch_size=10)
10
11
12 # Predict on the test set
13 y_pred_2_3 = model_2_3.predict(x_valid_k)
14 y_pred_classes = np.argmax(y_pred_2_3, axis=1)
15
16
17 cm = confusion_matrix(y_valid_k, y_pred_classes)
18
19 target_names = ['Normal', 'IR007', 'BB007', 'OR007']
20
21 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=
target_names)
22 disp.plot(cmap=plt.cm.Blues)
23 plt.title('Confusion Matrix')
24 plt.show()
25
26
27 accuracy = accuracy_score(y_valid_k, y_pred_classes)
28
29 accuracies.append(accuracy)
30
31 # Calculate the average accuracy
32 average_accuracy = np.mean(accuracies)
33 print(f'Average Accuracy: {average_accuracy:.2f}')
```

Code 27: import dataset and library

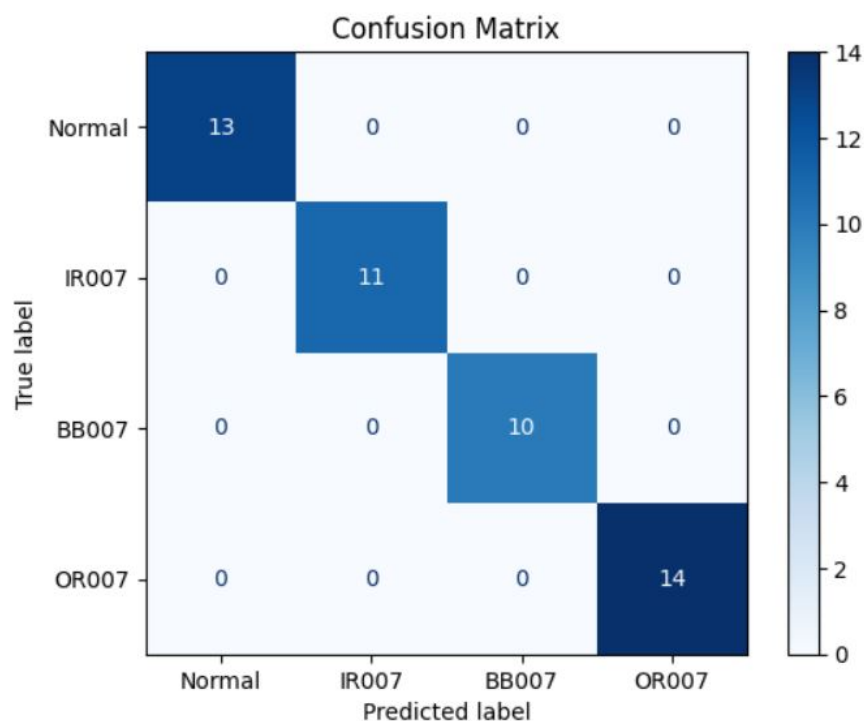


نتایج ماتریس‌های درهم‌ریختگی به صورت شکل ۲۱ قابل مشاهده است.



شکل ۲۱: نتیجه آموزش و ارزیابی با استفاده از SGD

در نهایت نتیجه ارزیابی روی مجموعه داده ارزیابی نیز به صورت شکل ۲۲ است.



شکل ۲۲: نتیجه آموزش و ارزیابی با استفاده از SGD

### ۳ سوال سوم

ابتدا کتابخانه‌های مورد نیاز را فراخوانی می‌کنیم. سپس فایل CSV داندودی از مرجع مورد ذکر سوال را با دستور gdown فراخوانی می‌کنیم و آن را بصورت یک دیتافریم درمی‌آوریم. دیتافریم تشکیل شده بصورت زیر است. هم‌طور که قابل مشاهده است دیتاست دارای ۵ ویژگی است که برخی به صورت عددی و برخی به صورت توصیفی هستند. همچنین تعداد کل داده‌ها برابر ۲۰۰ است.

```
1 !gdown 1XI-a6XbUk23bh0dNirGN8-drYX1ZUtFj
2 drug_data = pd.read_csv("/content/drug200.csv")
3
4 drug_data.head()
```

Code 28: import dataset and library

در صورت سوال ذکر شده است که بخشی از داده‌ها مورد استفاده قرار گیرد. بنابراین لازم است ابتدا اطلاعاتی در مورد داده‌ها به دست آوریم تا بتوانیم در صورت لزوم این کار را انجام دهیم. ابتدا مقدار داده‌های موجود در هر کلاس را مورد بررسی قرار می‌دهیم. کد و خروجی به صورت زیر است.

```
1 column_name = 'Drug'
2 num_unique_values = drug_data[column_name].nunique()
```



	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY

شکل ۲۳: نتیجه آموزش و ارزیابی با استفاده از SGD

```

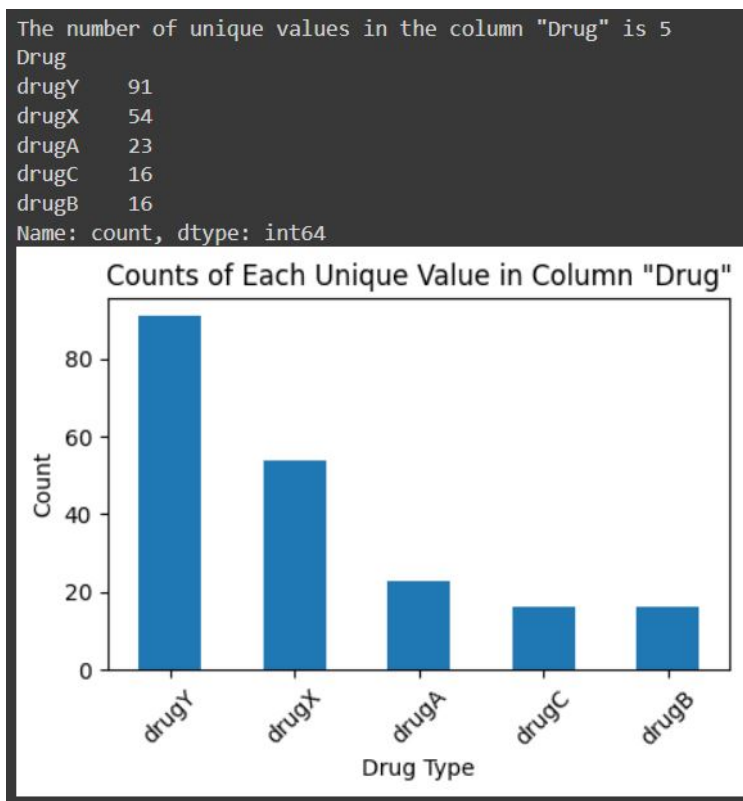
3 value_counts = drug_data[column_name].value_counts()
4 print(f'The number of unique values in the column "{column_name}" is {
    num_unique_values}')
5 print(value_counts)
6
7 # Plotting the value counts
8 plt.figure(figsize=(10, 6))
9 value_counts.plot(kind='bar')
10 plt.title(f'Counts of Each Unique Value in Column "{column_name}"')
11 plt.xlabel('Drug Type')
12 plt.ylabel('Count')
13 plt.xticks(rotation=45)
14 plt.show()

```

Code 29: import dataset and library

تعداد داده‌ها موجود در دیتاست متناسب با یکدیگر نیست. بیشترین مقدار داده در یک طبقه برابر ۹۱ و کمترین برابر ۱۶ است. ممکن است این عدم تناسب برای شبکه مشکل ایجاد کند. البته می‌توان ابتدا پیش از متناسب کردن تعداد داده‌ها با یکدیگر پس از آموزش، مدل را مورد ارزیابی قرار داد و در صورت لزوم این چنین پیش پردازشی روی مدل انجام دهیم. پیش از تقسیم داده‌ها نیاز است که ویژگی‌های دیتاست را مطابق با مدل درخت تصمیم آماده کنیم. لازم است که نوع دیتای ویژگی‌های





شکل ۲۴: نتیجه آموزش و ارزیابی با استفاده از SGD

دیتاست از یک نوع باشد. بنابراین ویژگی‌هایی که از جنس توصیفی هستند را به صورتی مجازی تغییر می‌دهیم و به هر کلاس یک عدد نسبت می‌دهیم. ویژگی‌های Sex, BP, Cholesterol توصیفی هستند. بنابراین ابتدا کلاس‌های هر کدام از ویژگی‌ها را لازم است بدانیم. کد زیر نحوه انجام این فرآیند و خروجی مورد نظر را نشان می‌دهد.

```
1 unique_values_Sex = np.unique(drug_data["Sex"])
2 unique_values_BP = np.unique(drug_data["BP"])
3 unique_values_Cholesterol = np.unique(drug_data["Cholesterol"])
4
5 print("unique_values_Sex = ", unique_values_Sex)
6 print("unique_values_BP = ", unique_values_BP)
7 print("unique_values_Cholesterol", unique_values_Cholesterol)
8
9
10 unique_values_Sex = ['F' 'M']
11 unique_values_BP = ['HIGH' 'LOW' 'NORMAL']
```



```
12 unique_values_Cholesterol ['HIGH' 'NORMAL']
```

Code 30: import dataset and library

سپس با توجه به کلاس های هر ویژگی توصیفی، به هر نوع کلاس یک عدد به خصوص نسبت می دهیم که مناسب دادن به مدل درخت تصمیم باشد. خروجی این تغییر در داده ها نیز به صورت **شکل ۲۵** است.

```
1 Sex_mapping = {'F': 1, 'M': 2}
2 BP_mapping = {'HIGH': 1, 'LOW': 2, 'NORMAL': 3}
3 Cholesterol_mapping = {'HIGH': 1, 'NORMAL': 2}
4
5 drug_data['Sex'] = drug_data['Sex'].replace(Sex_mapping)
6 drug_data['BP'] = drug_data['BP'].replace(BP_mapping)
7 drug_data['Cholesterol'] = drug_data['Cholesterol'].replace(
    Cholesterol_mapping)
8
9 drug_data.head()
```

Code 31: import dataset and library

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	1	1	1	25.355	drugY
1	47	2	2	1	13.093	drugC
2	47	2	2	1	10.114	drugC
3	28	1	3	1	7.798	drugX
4	61	1	2	1	18.043	drugY

شکل ۲۵: نتیجه آموزش و ارزیابی با استفاده از SGD

در قدم بعدی باید از مجموعه داده مورد نظر مقادیر مربوط به ویژگی ها را در قالب  $X$  و برچسب ها در قالب  $y$  درآورده شود. سپس داده های مربوط به آموزش و ارزیابی از یکدیگر جدا شده اند. نسبت داده های ارزیابی ۱۵ درصد در نظر گرفته شده است.

```
1 X_drug = drug_data.drop(columns=['Drug']).values
2 y_drug = drug_data['Drug'].values
```



```

3
4 X_train_d, X_test_d, y_train_d, y_test_d = train_test_split(X_drug, y_drug,
    random_state=4, test_size=0.15)
5 X_train_d.shape, X_test_d.shape
6
7 ((170, 5), (30, 5))

```

Code 32: import dataset and library

گام بعدی طراحی درخت تصمیمی است که کلاس‌های مورد نظر داده‌های ما را از یکدیگر جدا کند. این درخت تصمیم را با استفاده از دستور `DecisionTreeClassifier` فراخوانی کرده‌ایم. در این دستور لازم است برخی پارامترها تعریف شوند. یکی از پارامترها `criterion` است که تابعی است که میزان کیفیت تقسیم‌بندی را اندازه‌گیری می‌کند. این مقادیر برای این پارامتر قابل استفاده هستند: "gini"، "entropy"، "log-loss" همچنین پارامتر `ccp-alpha` نیز یک پارامتر مهم است. این پارامتر در مورد میزان پیچیدگی مورد استفاده برای هرس حداقلی `Cost-Complexity` است. به این صورت که در فرآیند آموزش، درختی فرعی با بیشترین میزان `Cost-Complexity` که از مقدار `ccp-alpha` کمتر باشد، انتخاب خواهد شد. مقادیر مورد استفاده در این مرحله در کد آمده است.

```

1 clf = tree.DecisionTreeClassifier(random_state=4, ccp_alpha=0.1, criterion='
    entropy')
2 clf.fit(X_train_d, y_train_d)
3 tree.plot_tree(clf)

```

Code 33: import dataset and library

```

DecisionTreeClassifier
DecisionTreeClassifier(ccp_alpha=0.5, criterion='entropy', random_state=4)

```

شکل ۲۶: نتیجه آموزش و ارزیابی با استفاده از SGD

درخت تصمیم طراحی شده با توجه به پارامترهای انتخابی به صورت **شکل ۲۷** است. در این حالت ۲ مرحله تصمیم‌گیری در درخت تصمیم صورت گرفته است. در مرحله اول با آستانه گذاری برای ویژگی `Na-to-K` انجام شده است. به این صورت که هر داده‌ای که دارای این ویژگی بالاتر از آستانه باشد متعلق به کلاس پنجم است و اگر کمتر از این آستانه باشد متعلق به سایر کلاس‌هاست. در نقطه تصمیم‌گیری بعدی یک آستانه گذاری برای ویژگی `BP` گذاشته شده است. که طبق آن اگر مقدار این ویژگی کمتر از آستانه باشد آن داده متعلق به کلاس اول یا دوم و اگر بیشتر از آستانه باشد متعلق به کلاس سوم یا چهارم است.

برای تست روی داده‌های ارزیابی نیز به صورت ۴ داریم که برابر مقدار بالایی نیست. برای حل این مشکل و اینکه مدل بتواند بصورت دقیق‌تری کلاس‌ها را جداسازی کند و مطابق خواسته سوال، پارامترهای طراحی درخت را تغییر می‌دهیم.

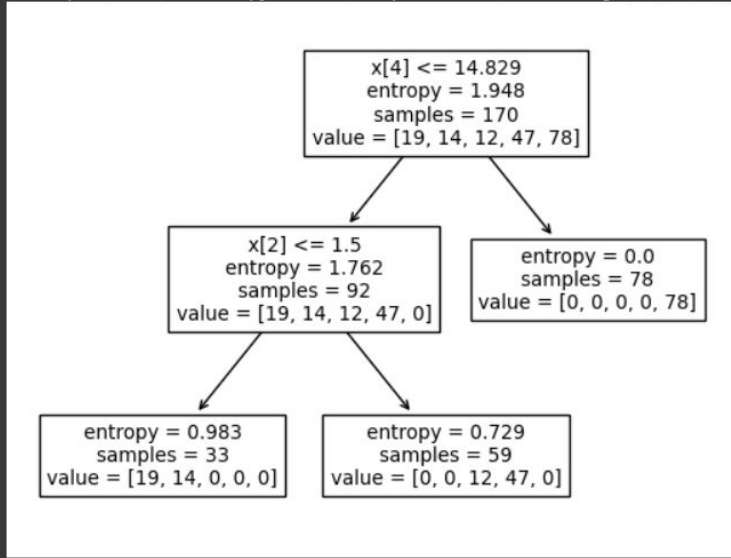
```

1 clf.score(X_test_d, y_test_d)

```



```
[Text(0.6, 0.8333333333333334, 'x[4] <= 14.829\nentropy = 1.948\nsamples = 170\nvalue = [19, 14, 12, 47, 78]'),
Text(0.4, 0.5, 'x[2] <= 1.5\nentropy = 1.762\nsamples = 92\nvalue = [19, 14, 12, 47, 0]'),
Text(0.2, 0.16666666666666666, 'entropy = 0.983\nsamples = 33\nvalue = [19, 14, 0, 0, 0]'),
Text(0.6, 0.16666666666666666, 'entropy = 0.729\nsamples = 59\nvalue = [0, 0, 12, 47, 0]'),
Text(0.8, 0.5, 'entropy = 0.0\nsamples = 78\nvalue = [0, 0, 0, 0, 78]')]
```



شکل ۲۷: نتیجه آموزش و ارزیابی با استفاده از SGD

2

3 0.8

Code 34: import dataset and library

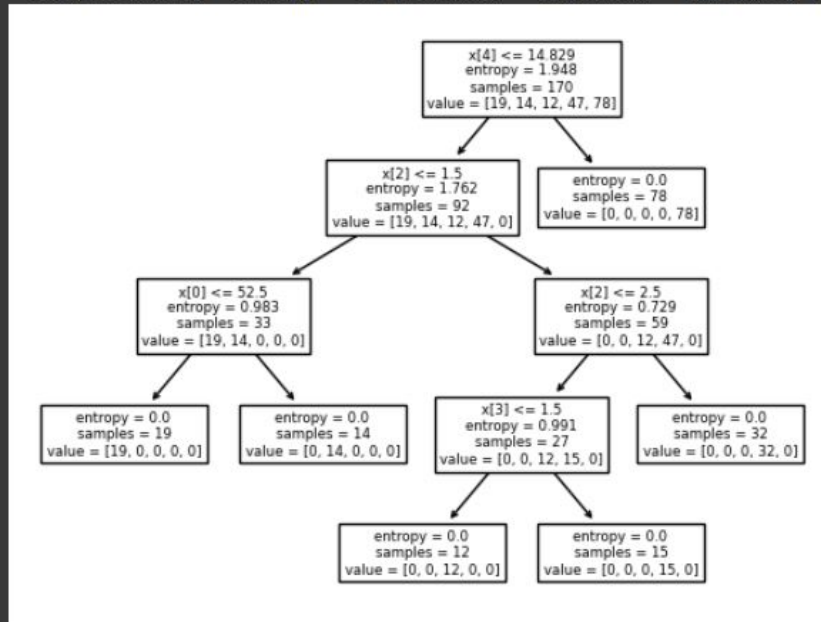
تغییر پارامتر را بصورت انتخاب مقدار ccp-alpha برابر ۱۰۰ انتخاب می‌دهیم. درخت طراحی شده پس از انتخاب این پارامتر بصورت **شکل ۲۸** درمی‌آید. با این طراحی با تعداد ۵ نقطه تصمیم‌گیری و روی چهار ویژگی مختلف، مدل توانسته است مطابق تصویر در تمامی لیف‌نود ها، تنها یک کلاس حضور داشته باشد که نشان‌دهنده تفکیک شدن همه کلاس‌ها از یکدیگر است.

۱.۳ با استفاده از ماتریس درهم ریختگی و حداقل سه شاخصه ارزیابی مربوط به وظیفه طبقه بندی، عمل کرد درخت آموزش داده شده خود را روی بخش آزمون داده‌ها ارزیابی کنید و نتایج را به صورت دقیق گزارش کنید. تأثیر مقادیر کوچک و بزرگ حداقل دو فرایارامتر را بررسی کنید. تغییر فرایارامترهای مربوط به هرس کردن چه تأثیری روی نتایج دارد و مزیت آن چیست؟

ماتریس در هم ریختگی و سه معیار recall, precision, accuracy برای ارزیابی عملکرد مدل روی بخش داده‌های تست مورد استفاده قرار گرفته‌اند. همانطور که از ماتریس درهم‌ریختگی در **شکل ۲۹** مشخص است تنها یک داده به اشتباه تشخیص داده شده است. این داده متعلق به کلاس B بوده و به اشتباه به کلاس A تعلق پیدا کرده است. این تشخیص اشتباه در معیارهای ذکر شده به صورت درصد خود را نشان داده‌اند.



```
[Text(0.625, 0.9, 'x[4] <= 14.829\nentropy = 1.948\nsamples = 170\nvalue = [19, 14, 12, 47, 78]'),
Text(0.5, 0.7, 'x[2] <= 1.5\nentropy = 1.762\nsamples = 92\nvalue = [19, 14, 12, 47, 0]'),
Text(0.25, 0.5, 'x[0] <= 52.5\nentropy = 0.983\nsamples = 33\nvalue = [19, 14, 0, 0, 0]'),
Text(0.125, 0.3, 'entropy = 0.0\nsamples = 19\nvalue = [19, 0, 0, 0, 0]'),
Text(0.375, 0.3, 'entropy = 0.0\nsamples = 14\nvalue = [0, 14, 0, 0, 0]'),
Text(0.75, 0.5, 'x[2] <= 2.5\nentropy = 0.729\nsamples = 59\nvalue = [0, 0, 12, 47, 0]'),
Text(0.625, 0.3, 'x[3] <= 1.5\nentropy = 0.991\nsamples = 27\nvalue = [0, 0, 12, 15, 0]'),
Text(0.5, 0.1, 'entropy = 0.0\nsamples = 12\nvalue = [0, 0, 12, 0, 0]'),
Text(0.75, 0.1, 'entropy = 0.0\nsamples = 15\nvalue = [0, 0, 0, 15, 0]'),
Text(0.875, 0.3, 'entropy = 0.0\nsamples = 32\nvalue = [0, 0, 0, 32, 0]'),
Text(0.75, 0.7, 'entropy = 0.0\nsamples = 78\nvalue = [0, 0, 0, 0, 78]')]
```



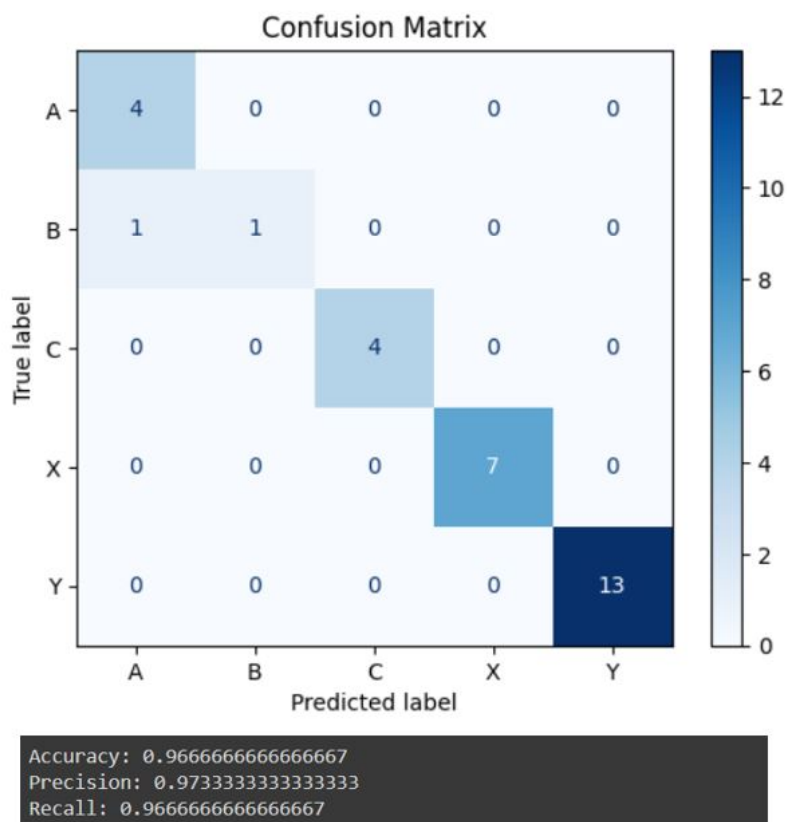
شکل ۲۸: نتیجه آموزش و ارزیابی با استفاده از SGD

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
4
5 y_pred_classes = np.argmax(y_pred_2_2, axis=1) # Convert predicted
    probabilities to class labels
6 cm = confusion_matrix(y_test_d, y_pred)
7 target_names = ['A', 'B', 'C', 'X', 'Y']
8
9 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=target_names
    )
10 disp.plot(cmap=plt.cm.Blues)
11 plt.title('Confusion Matrix')
```



```
12 plt.show()
13
14
15 from sklearn.metrics import accuracy_score, precision_score, recall_score
16
17 accuracy = accuracy_score(y_test_d, y_pred)
18 precision = precision_score(y_test_d, y_pred, average='weighted')
19 recall = recall_score(y_test_d, y_pred, average='weighted')
20 print("Accuracy:", accuracy)
21 print("Precision:", precision)
22 print("Recall:", recall)
```

Code 35: import dataset and library



شکل ۲۹: نتیجه آموزش و ارزیابی با استفاده از SGD

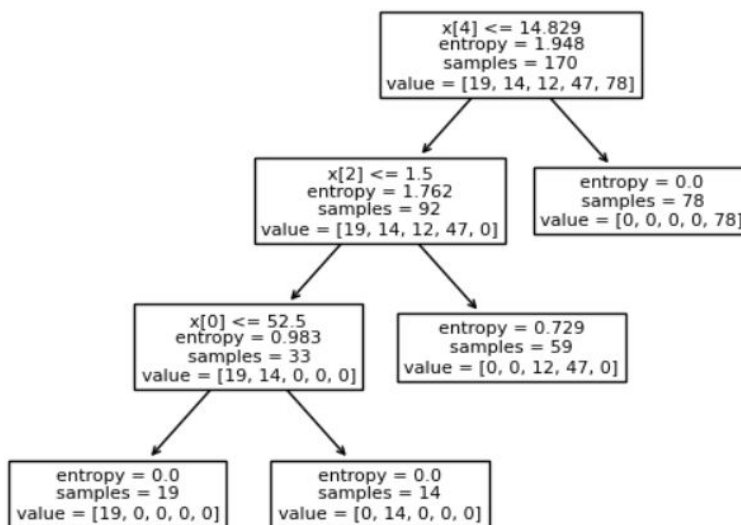
اثر تغییر یکی از پارامترها که ccp-alpha بود مورد بررسی قرار گرفت. این پارامتر پیچیدگی درخت تصمیم را از طریق هرس حداقل Cost-Complexity کنترل می کند. مقدار کوچکتر ccp-alpha به این معنی است که جریمه اضافه کردن گره به درخت کمتر است.



بنابراین، الگوریتم به احتمال زیاد گره‌ها را اضافه می‌کند و درخت بزرگ‌تری را رشد می‌دهد، که به طور بالقوه منجر به مدل پیچیده‌تری می‌شود که می‌تواند بهتر با داده‌های آموزشی تناسب داشته باشد. با این حال، این افزایش پیچیدگی همچنین ممکن است منجر به تطبیق بیش از حد داده‌های آموزشی و کاهش عملکرد در داده‌های دیده نشده شود.

به عنوان پارامتر دوم max-depth را انتخاب می‌کنیم. این پارامتر تعیین‌کننده عمق درخت تولید شده است. هر چقدر مقدار این پارامتر بیشتر باشد، درخت می‌تواند تعداد نودهای تصمیم‌گیری بیشتری داشته باشد. معمولاً از این پارامتر برای جلوگیری از اورفیت زدن استفاده می‌شود که درخت خیلی روی دیتاهای ارزیابی تطبیق بیش از اندازه نیابد. در این سوال چون ما با عمق ۴ توانستیم به عملکرد بسیار مطلوبی از مدل برسیم، این پارامتر را برابر ۳ قرار می‌دهیم. البته این کار صرفاً برای مشاهده نتیجه این تغییر پارامتر صورت می‌گیرد و ما در قسمت قبل به نتیجه دلخواه مدل خود رسیده‌ایم.

با اعمال max-depth برابر ۳ این درخت تصمیم را داریم:



شکل ۳۰: نتیجه آموزش و ارزیابی با استفاده از SGD

مشاهده می‌شود که با اعمال این پارامتر تعداد نودهای تصمیم‌گیرنده کاهش یافته و لیف‌نودهایی با کلاس‌های مشترک دیده می‌شود.

۲.۳ توضیح دهید که روش‌هایی مانند جنگل تصادفی و AdaBoost چگونه می‌توانند به بهبود نتایج کمک کنند. سپس، با انتخاب یکی از این روش‌ها و استفاده از فرآیندهای مناسب، سعی کنید نتایج پیاده‌سازی در مراحل قبلی را ارتقاء دهید.

هر دو روش از تکنیک‌های یادگیری جمعی هستند که در موضوع درخت تصمیم برای رسیدن به نتایج بهبودیافته‌تر مورد استفاده قرار می‌گیرند.

**Random Forest**: در این روش چندین درخت تصمیم در طول آموزش ساخته می‌شود. هر درخت به صورت مستقل و در حالی که فرآیند انتخاب ویژگی به صورت تصادفی است، ساخته می‌شود. هر درخت به صورت مستقل یک پیش‌بینی انجام می‌دهد. پیش‌بینی نهایی در تسک طبقه‌بندی به صورت رای‌گیری بیشینه صورت می‌پذیرد. این تکنیک به دلیل استحکام در برابر اورفیت شدن و توانایی در



برابر مدیریت داده‌های بزرگ شناخته می‌شود.

**AdaBoost**: این تکنیک با ترکیب چند درخت تصمیم ضعیف برای ایجاد یک درخت تصمیم قوی کار می‌کند. این روش با تنظیم وزن داده‌های آموزشی بر اساس نحوه عملکرد درخت تصمیم پیشین کار می‌کند. در درخت ابتدایی اولیه به هر داده وزن یکسانی داده می‌شود. در تکرارهای بعدی تمرکز و وزن بیشتری برای داده‌های با طبقه‌بندی اشتباه گذاشته می‌شود و به مدل اجازه می‌دهد از اشتباهات قبلی خود درس بگیرد. پیش‌بینی نهایی با ترکیب طبقه‌بندهای قبلی با اختصاص دادن وزن مخصوص به هر یک شکل می‌گیرد. از تکنیک Random Forest استفاده می‌کنیم. با استفاده از روش GridSearch پارامترهای مختلف را در این روش مورد ارزیابی قرار می‌دهیم.

```
1 param_grid = {
2     'n_estimators': [50, 100, 150],
3     'bootstrap': [True, False]
4     'max_depth': [6, 10, 20, None],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 4],
7 }
8
9 rf_classifier = RandomForestClassifier(random_state=4)
10 grid_search = GridSearchCV(estimator=rf_classifier, param_grid=param_grid, cv
    =5, scoring='accuracy')
11 grid_search.fit(X_train_d, y_train_d)
12
13 # Get the best hyperparameters
14 best_params = grid_search.best_params_
15 print("Best Hyperparameters:", best_params)
16
17 # Use the best model to predict on test data
18 best_model = grid_search.best_estimator_
19 y_pred = best_model.predict(X_test_d)
20
21 # Evaluate the classifier
22 accuracy = accuracy_score(y_test_d, y_pred)
23 print("Accuracy:", accuracy)
```

Code 36: import dataset and library

نتایج نهایی به صورت شکل ۳۱ قابل مشاهده است.





```
Best Hyperparameters: {'bootstrap': True, 'max_depth': 6, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}
Accuracy: 0.9333333333333333
```

شکل ۳۱: نتیجه آموزش و ارزیابی با استفاده از SGD

## ۴ عنوان سوال چهارم

ابتدا کتابخانه‌های مورد نیاز را فراخوانی می‌کنیم.

```
1 from pandas import DataFrame as df
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 %matplotlib inline
```

Code 37: import libraries (Python)

سپس دیتاست بیماری قلب با استفاده از دستور gdown قابلیت دسترسی پیدا می‌کند. سپس با کتابخانه pandas این دیتاست خوانده می‌شود. با دستور head() عنوان ستون‌های دیتاست قابل مشاهده خواهند بود.

```
1 !gdown 1qYg0L0-iDvhawJw-v5keUtQc9FpYALxB
2
3 data = pd.read_csv("/content/heart.csv")
4 data.head()
```

Code 38: import libraries (Python)

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3	0
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	0
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3	0
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	0

شکل ۳۲: بلوک دیاگرام مدل طبقه‌بند خطی

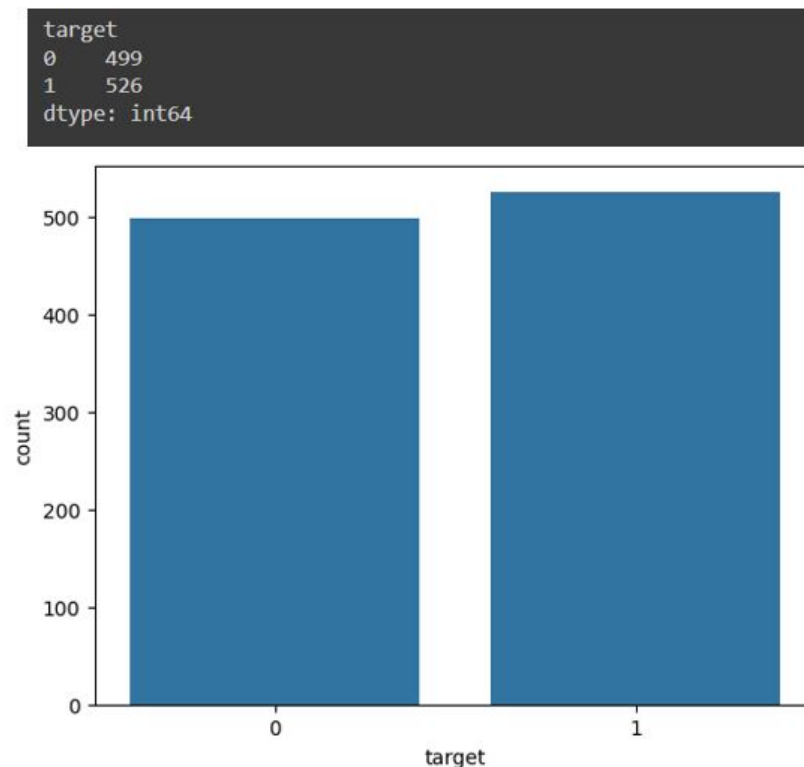
با استفاده از این دستورات می‌توان به میزان داده‌های موجود در هر کلاس پی برد.



```
1 data.groupby('target').size()  
2 sns.countplot(x='target', data=data)
```

Code 39: import libraries (Python)

خروجی این کد بصورت زیر است. همانطور که قابل مشاهده است میزان داده‌ها در دو کلاس نامتوازن است. بنابراین می‌توان در نظر داشت که ممکن است کار متوازن کردن داده‌ها در نتایج تاثیر داشته باشد.



شکل ۳۳: بلوک دیاگرام مدل طبقه‌بند خطی

با استفاده از این دستورها می‌توان اطلاعات ارزشمندی از دیتاست به دست آورد.

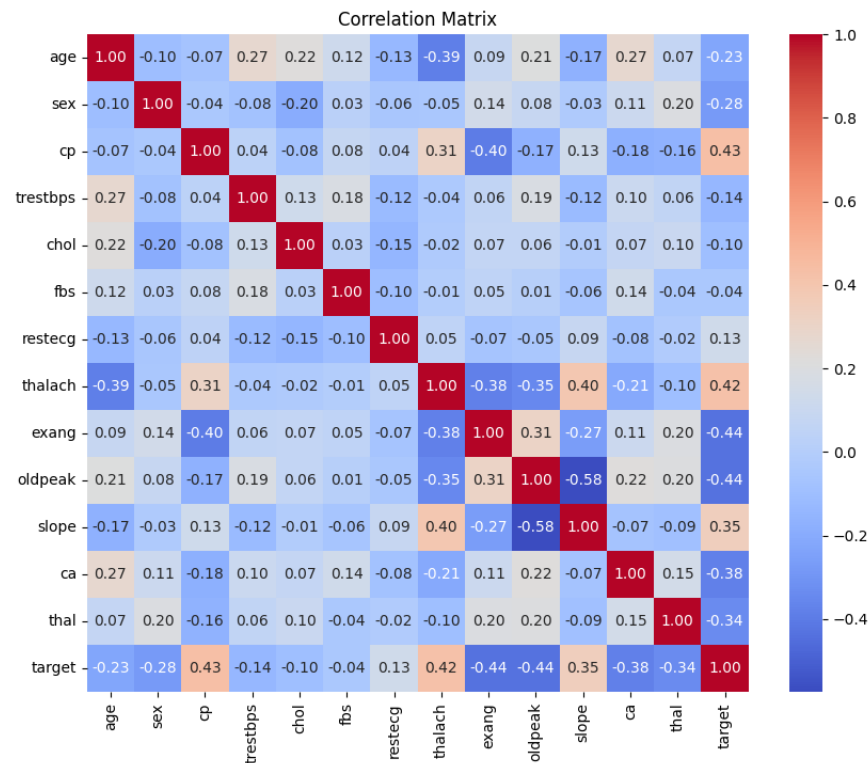
```
1 sns.pairplot(data, hue='target')  
2  
3 data.hist()  
4  
5 correlation_matrix = data.corr()  
6 plt.figure(figsize=(10, 8))  
7 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")  
8 plt.title('Correlation Matrix')
```



```
plt.show()
```

Code 40: import libraries (Python)

خط اول شبکه‌ای از نمودارهای پراکنده ایجاد می‌کند که در آن هر نمودار نشان‌دهنده رابطه بین جفت‌های متغیر در مجموعه داده است و نقاط بر اساس مقادیر در ستون 'هدف' رنگ می‌شوند. این تجسم به ویژه برای بررسی روابط بین چندین متغیر و نحوه ارتباط آنها با متغیر هدف مفید است. خط دوم یک هیستوگرام جداگانه برای هر ستون عددی ایجاد می‌کند و این امکان را می‌دهد که توزیع مقادیر را برای هر متغیر به صورت بصری بررسی کرد. قسمت سوم این کد یک نمایش بصری از همبستگی‌های بین متغیرهای مختلف در مجموعه داده ایجاد می‌کند و به شناسایی الگوها و روابط بین متغیرها کمک می‌کند. ماتریس همبستگی به صورت زیر قابل مشاهده است.



شکل ۳۴: بلوک دیاگرام مدل طبقه‌بند خطی

همانطور که مشاهده می‌شود ارتباط بین ویژگی `thalach` و `slope` نسبت به سایر ویژگی‌ها بیشتر است و می‌توان از این ویژگی در ادامه برای کارهای پیش‌پردازشی بهره برد. یکی از کارهایی که به طور معمول در پیش‌پردازش استفاده می‌شود، شافل کردن داده‌هاست. ممکن است ترتیب چیده شدن دیتاها پشت یکدیگر از قبل مناسب نباشد و مدل را دچار اشتباه کند. شافل کردن باعث می‌شود مدل به دیتاهای مربوط به یک کلاس واحد وابسته نشود.

```
from sklearn.utils import shuffle
```



```
2 shuffled_data = shuffle(data, random_state=4)
3 print("shuffled_data shape = ",shuffled_data.shape)
4
5 shuffled_data shape = (1025, 14)
```

Code 41: import libraries (Python)

سپس نیاز است که مجموعه داده را به دو بخش ویژگی و تارگت تقسیم کنیم. تارگت در ستونی با عنوان target قابل دسترسی است.

```
1 X=np.array(shuffled_data.loc[:,data.columns!='target'])
2 y=np.array(shuffled_data.loc[:,data.columns=='target'])
3 print(X.shape, y.shape)
4
5 (1025, 13) (1025, 1)
```

Code 42: import libraries (Python)

با توجه به نامتوازن بودن تعداد داده‌های متعلق به هر کلاس، می‌توان این دو کلاس را با روش Random Undersampling متوازن کرد. به این صورت که از کلاس با داده‌های بیشتر بصورت رندوم حذف داده صورت بگیرد تا تعداد داده‌های هر کلاس برابر شود. البته با توجه به اینکه این میزان تفاوت داده‌های هر کلاس چشم‌گیر نیست، احتمالاً اثر چندانی در نتیجه نخواهد داشت. البته کد طوری در نظر گرفته شده است که بتوان هم بدون توازن و هم پس از توازن داده کار طبقه‌بندی را انجام داد.

```
1 from imblearn.under_sampling import RandomUnderSampler
2
3 # rus = RandomUnderSampler(random_state=4)
4 rus = RandomUnderSampler(sampling_strategy={0: 499, 1: 499}, random_state=4)
5
6 # Resample the dataset
7 UnderSample = True
8 if UnderSample :
9     X_resampled, y_resampled = rus.fit_resample(X, y)
10    print("Class Distribution after Random Undersampling:")
11    print(pd.Series(y_resampled).value_counts())
12 else :
13     X_resampled = X
14     y_resampled = y
```

Code 43: import libraries (Python)



خروجی پس از توازن داده‌ها بصورت دوش Random Undersampling بصورت زیر است.

```
Class Distribution after Random Undersampling:
0      499
1      499
Name: count, dtype: int64
```

شکل ۳۵: بلوک دیاگرام مدل طبقه‌بند خطی

تقسیم داده‌ها به دو بخش آموزش و آزمون بصورت زیر قابل انجام است. با توجه به تعداد حدود ۵۰۰ داده، مناسب است که مقدار ۲۰ درصد از داده‌ها را به مجموعه داده تست اختصاص دهیم.

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled,
3     test_size=0.2, random_state=4)
4
5 print('train:', X_train.shape, y_train.shape, '\ntest: ', X_test.shape, y_test
6     .shape)
```

Code 44: import libraries (Python)

با توجه به اطلاعات مربوط به ویژگی‌های داده در شکل ۳۲ و با استفاده از دستور `data.describe()` می‌توان مشاهده کرد که بازه مقادیر ویژگی‌های مختلف تفاوت زیادی دارد و این موضوع می‌تواند چالش آفرین باشد و مناسب‌تر است که بازه ویژگی‌های مجموعه داده را با هم متناسب کرد. از روش `MinMaxScaler` بازه ویژگی‌های دیتاست را بین ۰ و ۱ می‌آوریم. البته باید توجه داشت که این کار صرفاً با اسفاده از داده‌های آموزش باشد و از اطلاعات داده‌های تست استفاده نشود.

```
1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2 scaler = MinMaxScaler()
3 scaler.fit(X_train)
4 # print(X_train[:5])
5 X_train_t = scaler.transform(X_train)
6 X_test_t = scaler.transform(X_test)
7 # print(X_train_t[:5])
```

Code 45: import libraries (Python)



در مرحله بعد به دو صورت کار طبقه‌بندی صورت می‌گیرد. با دستور آماده و استفاده از کتابخانه sklearn و همچنین به صورت کدزنی کلاس این طبقه‌بند. و فراخوانی و آموزش روی داده‌های آموزش آورده شده است.

```
1 class NaiveBayes:
2     def fit(self, X, y):
3         n_samples, n_features = X.shape
4         self._classes = np.unique(y)
5         n_classes = len(self._classes)
6
7         # calculate mean, var, and prior for each class
8         self._mean = np.zeros((n_classes, n_features), dtype=np.float64)
9         self._var = np.zeros((n_classes, n_features), dtype=np.float64)
10        self._priors = np.zeros(n_classes, dtype=np.float64)
11
12        for idx, c in enumerate(self._classes):
13            X_c = X[y == c]
14            self._mean[idx, :] = X_c.mean(axis=0)
15            self._var[idx, :] = X_c.var(axis=0)
16            self._priors[idx] = X_c.shape[0] / float(n_samples)
17
18
19        def predict(self, X):
20            y_pred = [self._predict(x) for x in X]
21            return np.array(y_pred)
22
23        def _predict(self, x):
24            posteriors = []
25
26            # calculate posterior probability for each class
27            for idx, c in enumerate(self._classes):
28                prior = np.log(self._priors[idx])
29                posterior = np.sum(np.log(self._pdf(idx, x)))
30                posterior = posterior + prior
31                posteriors.append(posterior)
32
33            # return class with the highest posterior
```



```

34         return self._classes[np.argmax posteriors)]
35
36     def _pdf(self, class_idx, x):
37         mean = self._mean[class_idx]
38         var = self._var[class_idx]
39         numerator = np.exp(-((x - mean) ** 2) / (2 * var))
40         denominator = np.sqrt(2 * np.pi * var)
41         return numerator / denominator
42
43
44 from sklearn.naive_bayes import GaussianNB
45 MyNB = NaiveBayes()
46 SKNB = GaussianNB()
47
48 MyNB.fit(X_train, y_train.ravel())
49 SKNB.fit(X_train, y_train.ravel())
50 pred = MyNB.predict(X_test)
51 pred2 = SKNB.predict(X_test)

```

Code 46: import libraries (Python)

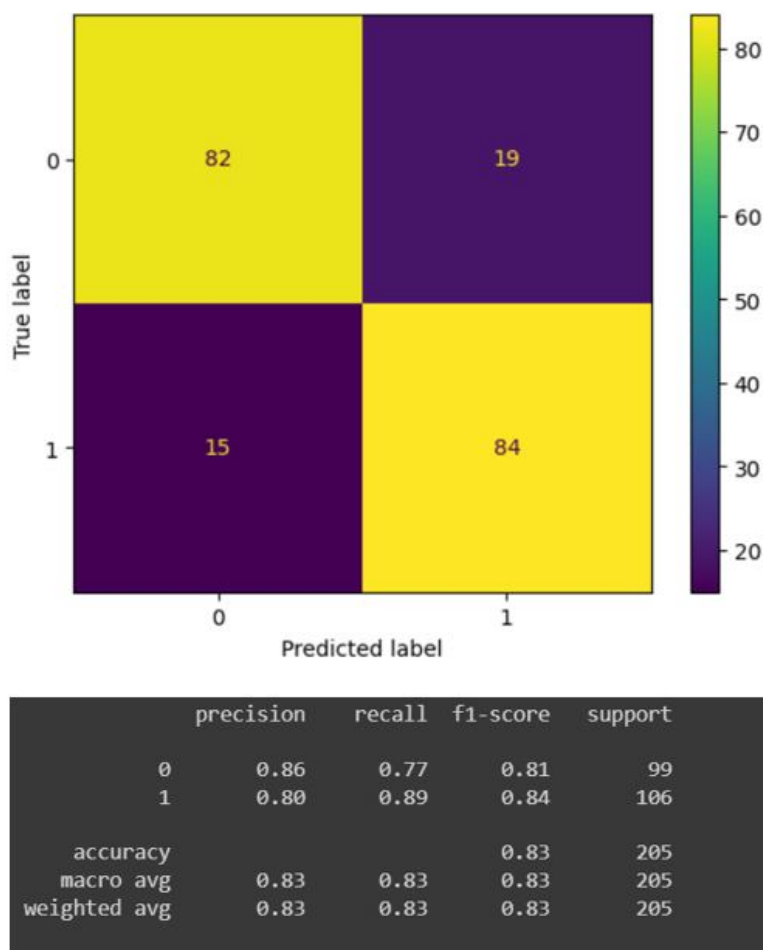
مطابق این کد پس از فراخوانی Bayes به دوروش گفته شده، با دستور fit این مدل را بر روی مجموعه داده آموزش ترین می کنیم. و پس از آن روی مجموعه تست کار ارزیابی را انجام می دهیم. سپس طبق خواسته مسئله ماتریس درهم ریختگی و سایر متریک ها را مور ارزیابی قرار می دهیم. مشاهده می شود که تعداد ۲۰۵ داده تست ارزیابی شده اند و مطابق معیارها مدل در تشخیص کلاس ۰ کمی عملکرد بهتری داشته است. در ماتریس درهم ریختگی نیز به ترتیب تعداد ۱۵ و ۱۹ داده متعلق به کلاس ۱ و ۰ به اشتباه به کلاس دیگری نسبت داده شده است و سایر داده ها به درستی طبقه بندی شده اند.

```

1 from sklearn.metrics import classification_report
2 print(classification_report(y_test, pred))
3
4 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
5 cm = confusion_matrix(y_test, pred)
6 names = list(data.groupby('target').groups.keys())
7 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=names)
8 disp.plot()
9 plt.show()

```

Code 47: import libraries (Python)



شکل ۳۶: بلوک دیاگرام مدل طبقه‌بند خطی

دو حالت Micro و Macro به روش‌های مختلف میانگین‌گیری معیارهای عملکرد برای مسائل طبقه‌بندی چند طبقه اشاره دارند. این روش‌های میانگین‌گیری برای جمع‌آوری نمرات دقت، فراخوانی و  $F1$  در چندین کلاس استفاده می‌شود. حالت Micro معیارها را بصورت گلوبال محاسبه می‌کند. این کار با شمارش تعداد TP، FN FP، بر روی همه کلاس‌ها انجام می‌پذیرد. در این حالت با همه نمونه‌ها به طور مساوی رفتار می‌شود، بنابراین به هر نمونه، صرف نظر از کلاس آن، وزن یکسانی داده می‌شود. در حالت Macro میانگین کلان معیارها را برای هر کلاس به طور مستقل محاسبه می‌شود و سپس میانگین این معیارها گرفته می‌شود. این حالت با همه کلاس‌ها به یک اندازه رفتار می‌کند، بنابراین به هر کلاس، صرف نظر از تعداد نمونه‌ها در هر کلاس، وزن یکسانی می‌دهد. حالت Micro زمانی مفید است که عملکرد کلی طبقه‌بندی کننده در همه موارد مورد نیاز باشد که ارزیابی شود. این به شدت تحت تأثیر عملکرد در کلاس‌های مختلف است. حالت Macro زمانی مفید است که عملکرد طبقه‌بندی کننده در هر کلاس به طور مستقل ارزیابی شود و با هر کلاس به طور مساوی رفتار شود. وقتی کلاس‌ها نامتعادل هستند، می‌تواند مناسب‌تر باشد. در انتها پنج داده به صورت تصادفی از مجموعه داده آزمون انتخاب کرده و خروجی پیش‌بینی شده را با خروجی واقعی مقایسه می‌کنیم. در این قسمت با دستور randint تعداد ۵ داده تولید می‌کنیم که بین ۰ تا طول بردار تعداد داده‌های تست باشد. از random-seed استفاده





کرده‌ایم که هر بار که اجرا می‌کنیم، همان داده‌های تولیدی اولیه را بدهد و هر بار مجدداً رندوم تولید نکند. خروجی نیز به صورت زیر قابل مشاهده است که تعداد ۴ تا از پیش‌بینی‌ها درست و تنها یکی از آنها متفاوت با لیبل خود دیتاست.

```
1 np.random.seed(4)
2 random_data = np.random.randint(0, len(X_test), 5)
3 # print(random_data)
4 five_data = X_test[random_data]
5 # five_data.shape
6
7 pred_t_five_data = MyNB.predict(five_data)
8 true_lable = y_test[random_data].T
9 print('pred: ', pred_t_five_data)
10 print('true lable: ', true_lable)
```

Code 48: import libraries (Python)

```
pred:      [0 1 1 0 0]
true lable: [[0 1 1 0 1]]
```

شکل ۳۷: بلوک دیاگرام مدل طبقه‌بند خطی

## منابع

[۱] صفحهٔ درس یادگیری ماشین.

[2] Steven X. Ding, “Data-driven Design of Fault Diagnosis and Fault-tolerant Control System”, Springer, 2014.

[3] S. Theodoridis and K. Koutroumbas, ”Pattern recognition”, Fourth Edition, Academic Press, 2009.