



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

۱	سوال اول	۸
۱.۱	فرآیند آموزش و ارزیابی یک مدل طبقه بند خطی را به صورت دیاگرامی بلوکی نمایش دهید و در مورد اجزای مختلف این دیاگرام بلوکی توضیحاتی بنویسید. تغییر نوع طبقه بندی از حالت دوکلاسه به چندکلاسه در کدام قسمت از این دیاگرام بلوکی تغییراتی ایجاد می کند؟ توضیح دهید.	۸
۲.۱	با استفاده از datasets.sklearn یک دیتاست با ۱۰۰۰ نمونه، ۴ کلاس و ۳ ویژگی تولید کنید و آن را به صورتی مناسب نمایش دهید. آیا دیتاستی که تولید کردید چالش برانگیز است؟ چرا؟ به چه طریقی می توانید دیتاست تولیدشده خود را چالش برانگیزتر و سخت تر کنید؟	۹
۳.۱	با استفاده از حداقل دو طبقه بند خطی آماده پایتون (در model-linear.sklearn) و در نظر گرفتن فرآیندهای مناسب، چهار کلاس موجود در دیتاست قسمت قبلی را از هم تفکیک کنید. ضمن توضیح روند انتخاب فرآیندها (مانند تعداد دوره آموزش و نرخ یادگیری)، نتیجه دقت آموزش و ارزیابی را نمایش دهید. برای بهبود نتیجه از چه تکنیک هایی استفاده کردید؟	۱۲
۴.۱	مرز و نواحی تصمیم گیری برآمده از مدل آموزش دیده خود را به همراه نمونه ها در یک نمودار نشان دهید. اگر می توانید نمونه هایی که اشتباه طبقه بندی شده اند را با شکل و رنگ متفاوت نمایش دهید.	۱۸
۵.۱	فرآیندی مشابه قسمت «۲» را با تعداد کلاس و ویژگی دلخواه؛ اما با استفاده از ابزار drawdata تکرار کنید. قسمت های «۳» و «۴» را برای این داده های جدید تکرار و نتایج را به صورتی مناسب نشان دهید.	۱۹
۲	عنوان سوال دوم	۲۳
۱.۲	با مراجعه به صفحه دیتاست CWRU Bearing با یک دیتاست مربوط به حوزه «تشخیص عیب» آشنا شوید. با جستجوی آن در اینترنت و مقالات، توضیحاتی از اهداف، ویژگی ها و حالت های مختلف این دیتاست ارائه کنید.	۲۳
۲.۲	قسمت دوم.	۲۳
۱.۲.۲	از هر کلاس نمونه با طول N جدا کنید (حداقل ۱۰۰ و N حداقل ۲۰۰ باشد). یک ماتریس از داده های هر دو کلاس به همراه برچسب مربوطه تشکیل دهید. می توانید پنجره ای به طول N در نظر بگیرید و در نهایت یک ماتریس $N \times$ از داده های هر کلاس استخراج کنید.	۲۳
۲.۲.۲	در مورد اهمیت استخراج ویژگی در یادگیری ماشین توضیحاتی بنویسید. سپس، با استفاده از حداقل ۸ عدد از روش های ذکر شده در جدول ۱، ویژگی های دیتاست قسمت «۲ آ» را استخراج کنید و یک دیتاست جدید تشکیل دهید.	۲۶
۳.۲.۲	ضمن توضیح اهمیت فرآیند ب رزدن (مخلوط کردن) داده ها را در صورت امکان مخلوط کرده و با نسبت تقسیم دلخواه و معقول به دو بخش «آموزش» و «ارزیابی» تقسیم کنید.	۲۹
۴.۲.۲	حداقل دو روش برای نرمال سازی داده ها را با ذکر اهمیت این فرآیند توضیح دهید و با استفاده از یکی از این روش ها، داده ها را نرمال کنید. آیا از اطلاعات بخش «ارزیابی» در فرآیند نرمال سازی استفاده کردید؟ چرا؟	۳۰



- ۳.۲ بدون استفاده از کتابخانه های آماده پایتون، مدل طبقه بند، تابع اتلاف و الگوریتم یادگیری و ارزیابی را کدنویسی کنید تا دو کلاس موجود در دیتاست به خوبی از یکدیگر تفکیک شوند. نمودار تابع اتلاف را رسم کنید و نتیجه ارزیابی روی داده های تست را با حداقل ۲ شاخصه محاسبه کنید. نمودار تابع اتلاف را تحلیل کنید. آیا می توان از روی نمودار تابع اتلاف و قبل از مرحله ارزیابی با قطعیت در مورد عمل کرد مدل نظر داد؟ چرا و اگر نمی توان، راه حل چیست؟ ۳۱
- ۴.۲ فرآیند آموزش و ارزیابی را با استفاده از یک طبقه بند خطی آماده پایتون (در `model-linear.sklearn`) انجام داده و نتایج را مقایسه کنید. در حالت استفاده از دستورات آماده سایکیت لرن، آیا راهی برای نمایش نمودار تابع اتلاف وجود دارد؟ پیاده سازی کنید. ۳۶
- ۵.۲ در مورد نرم افزار داده کاوی Orange و قابلیت های آن تحقیق کنید و سعی کنید این سوال یا یک مثال ساده تر را با استفاده از این نرم افزار پیاده سازی کنید (راهنمایی: می توانید از پیوندهای ۱ و ۲ کمک بگیرید). پاسخ به این قسمت از سوال اختیاری و امتیازی است. می توانید عملکرد خود را به صورت تصویری و یا ویدیویی هم نشان دهید. مقدار نمره امتیازی، وابسته به جامعیت مثال بررسی شده و استفاده از ویژگی های مختلف این ابزار است. ۳۸
- ۳ سوال سوم ۳۹
- ۱.۳ ابتدا هیت مپ ماتریس همبستگی و هیستوگرام پراکندگی ویژگی ها را رسم و تحلیل کنید. ۳۹
- ۲.۳ روی این دیتاست، تخمین LS و RLS را با تنظیم پارامترهای مناسب اعمال کنید. نتایج به دست آمده را با محاسبه خطاها و رسم نمودارهای مناسب برای هر دو مدل با هم مقایسه و تحلیل کنید. ۴۳
- ۳.۳ در مورد Weighted Least Square توضیح دهید و آن را روی دیتاست داده شده اعمال کنید. ۵۰
- ۴.۳ در مورد الگوریتم QR-Decomposition-Based RLS تحقیق کنید. پاسخ به این قسمت از سوال اختیاری و امتیازی است. ۵۲
- ۴ عنوان سوال چهارم ۵۳
- ۵ عنوان سوال پنجم ۵۳
- ۶ عنوان سوال ششم ۵۳
- ۱.۶ عنوان بخش اول سوال ششم ۵۳
- ۲.۶ عنوان بخش دوم سوال ششم ۵۳
- ۷ عنوان سوال هفتم ۵۴
- ۸ عنوان سوال هشتم ۵۵
- ۱.۸ عنوان بخش اول سوال هشتم ۵۵
- ۲.۸ عنوان بخش دوم سوال هشتم ۵۵
- ۳.۸ عنوان بخش سوم سوال هشتم ۵۵
- ۹ ضمیمه ۵۵



فهرست تصاویر

۸	بلوک دیاگرام مدل طبقه‌بند خطی	۱
۱۰	3D plot of dataset	۲
۱۲	2D plot of dataset	۳
۱۳	n-clusters-per-class=2, class-sep=10	۴
۱۴	بازه تغییرات ویژگی‌های دیتاست	۵
۱۵	نتیجه آموزش با پارامترهای مختلف	۶
۱۵	نتیجه آموزش با پارامترهای مختلف با class-sep=۲	۷
۱۵	نتایج آموزش و ارزیابی	۸
۱۶	نتایج آموزش با پارامترهای مختلف	۹
۱۶	نتایج آموزش و تست SGD	۱۰
۱۸	نتایج آموزش و تست SGD	۱۱
۱۹	دیتاست طراحی شده با drawdata	۱۲
۲۱	نتیجه آموزش و ارزیابی با استفاده از Regression Logistic	۱۳
۲۱	نتیجه آموزش با استفاده از SGD	۱۴
۲۲	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۵
۲۲	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۶
۲۴	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۷
۲۵	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۸
۲۵	نتیجه آموزش و ارزیابی با استفاده از SGD	۱۹
۲۷	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۰
۲۸	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۱
۲۹	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۲
۳۰	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۳
۳۰	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۴
۳۲	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۵
۳۳	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۶
۳۴	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۷
۳۴	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۸
۳۶	نتیجه آموزش و ارزیابی با استفاده از SGD	۲۹
۳۷	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۰
۳۸	نرم افزار Orange	۳۱
۴۰	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۲
۴۱	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۳
۴۲	نتیجه آموزش و ارزیابی با استفاده از SGD	۳۴



۳۵	نتیجه آموزش و ارزیابی با استفاده از SGD	۴۳
۳۶	نتیجه آموزش و ارزیابی با استفاده از SGD	۴۵
۳۷	نتیجه آموزش و ارزیابی با استفاده از SGD	۴۶
۳۸	نتیجه آموزش و ارزیابی با استفاده از SGD	۴۸
۳۹	نتیجه آموزش و ارزیابی با استفاده از SGD	۴۹
۴۰	نتیجه آموزش و ارزیابی با استفاده از SGD	۵۰
۴۱	نتیجه آموزش و ارزیابی با استفاده از SGD	۵۱
۴۲	نتیجه آموزش و ارزیابی با استفاده از SGD	۵۲
۴۳	شکل شماره ۱	۵۳



فهرست جداول

۱	جدول شماره ۱	۵۳
---	--------------	----



فهرست برنامه‌ها

۹ (Python) libraries import	۱
۹ Dataset Generate	۲
۱۰ ۳D in Dataset Plot	۳
۱۱ ۲D in Dataset Plot	۴
۱۲ dataset split	۵
۱۳ range feature	۶
۱۴ model train	۷
۱۴ score test and train	۸
۱۵ SGDClassifier	۹
۱۶ score test and train SGD	۱۰
۱۸ library drawdata Import	۱۱
۱۹ library drawdata Import	۱۲
۱۹ drawdata	۱۳
۲۰ dataframe creat	۱۴
۲۰ Datase Split	۱۵
۲۰ LR by Train	۱۶
۲۱ SGD by Train	۱۷
۲۲ SGD by Train	۱۸
۲۳ library and dataset import	۱۹
۲۴ library and dataset import	۲۰
۲۵ library and dataset import	۲۱
۲۵ library and dataset import	۲۲
۲۶ library and dataset import	۲۳
۲۷ library and dataset import	۲۴
۲۹ library and dataset import	۲۵
۲۹ library and dataset import	۲۶
۳۰ library and dataset import	۲۷
۳۱ library and dataset import	۲۸
۳۱ library and dataset import	۲۹
۳۲ library and dataset import	۳۰
۳۲ library and dataset import	۳۱
۳۴ library and dataset import	۳۲
۳۵ library and dataset import	۳۳
۳۶ library and dataset import	۳۴



۳۶	library and dataset import	۳۵
۳۹	library and dataset import	۳۶
۳۹	library and dataset import	۳۷
۳۹	library and dataset import	۳۸
۴۱	library and dataset import	۳۹
۴۱	library and dataset import	۴۰
۴۳	library and dataset import	۴۱
۴۴	library and dataset import	۴۲
۴۵	library and dataset import	۴۳
۴۶	library and dataset import	۴۴
۴۸	library and dataset import	۴۵
۵۰	library and dataset import	۴۶
۵۱	library and dataset import	۴۷
۵۴	(Python) Caption My	۴۸
۵۴	(MATLAB) Caption My	۴۹
۵۴	(C++) Caption My	۵۰
۵۴	(C) Caption My	۵۱



لینک پوشه گیت هاب:

https://github.com/ImanGandomi/MachineLearning2024/tree/main/ml_MiniProject_1

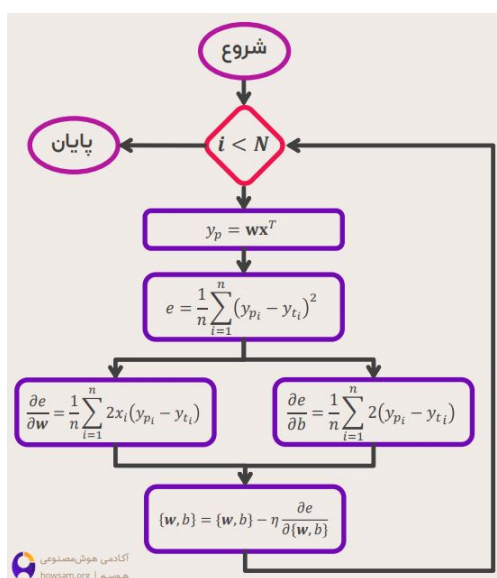
لینک گوگل کولب:

<https://colab.research.google.com/drive/1NwZ3BcL3-qJCNFXE9iJAdm-rj0CTM8ec?usp=sharing>

۱ سوال اول

۱.۱ فرآیند آموزش و ارزیابی یک مدل طبقه بند خطی را به صورت دیاگرامی بلوکی نمایش دهید و در مورد اجزای مختلف این دیاگرام بلوکی توضیحاتی بنویسید. تغییر نوع طبقه بندی از حالت دوکلاسه به چندکلاسه در کدام قسمت از این دیاگرام بلوکی تغییراتی ایجاد می کند؟ توضیح دهید.

ساده ترین بلوک دیاگرام یک طبقه بند خطی را می توان بصورت زیر در نظر گرفت. مطابق این سیستم در محله اول با تعریف مجموعه داده مورد نظر و تعداد دوره آموزشی، مقداردهی اولیه پارامترها صورت می پذیرد. این پارامترها شامل مقادیر وزن ها و بایاس ها است. در قدم بعدی مقدار وزن و ورودی داده شده از مجموعه داده مورد ارزیابی و محاسبه قرار گرفته و مقدار خروجی $y-h$ محاسبه می شود. در قدم بعدی نوبت به محاسبه اتلاف می رسد. اتلاف در واقع اختلاف مقدار پیش بینی شده و مقدار واقعی برچسب داده مورد نظر است. این اتلاف با توجه به تابع اتلاف مورد نظر ما محاسبه می شود. سپس باید گرادیان پارامترها محاسبه شود. در واقع گرادیان مقدار مشتق آن پارامترها است. همانطور که در ریاضیات برای رسیدن به مقدار بهینه یه پارامتر از آن مشتق می گرفتیم، اینجا هم همین اتفاق می افتد. در ادامه این مقدار محاسبه شده برای گرادیان این پارامترها باید با در نظر گرفتن یک مقدار کنترل شده ای به اسم نرخ یادگیری، از مقدار قبلی آن پارامتر کاسته شود تا به سمت مقادیر بهینه حرکت کنیم. این کار تا پایان دوره های آموزشی تکرار می شود. مرحله ارزیابی پس از متوقف شدن و ذخیره شدن بهترین مدل انجام می پذیرد. به این صورت که وزن های بهترین مدل که ذخیره شده اند، مورد استفاده قرار گرفته و روی داده های ارزیابی تست می شود.



شکل ۱: بلوک دیاگرام مدل طبقه بند خطی



اگر بخواهیم در یکی مدل طبقه‌بند خطی نوع طبقه‌بندی را از حالت دوکلاسه به حالت چندکلاسه تغییر دهیم، در قسمت تعریف تابع اتلاف باید تغییراتی ایجاد کنیم. در تصویر مربوطه از تابع اتلاف BCE استفاده شده است که مناسب حالت‌های دوکلاسه است. این قسمت باید با توجه به چندکلاسه بودن تغییر کند و تابعی مورد استفاده قرار گیرد که مناسب باشد. همچنین هر قسمتی که دارای تعریف دوکلاسه باشد باید تغییر کرده و حالت چندکلاسه به خود بگیرد.

۲.۱ با استفاده از datasets.sklearn یک دیتاست با ۱۰۰۰ نمونه، ۴ کلاس و ۳ ویژگی تولید کنید و آن را به صورتی مناسب نمایش دهید. آیا دیتاستی که تولید کردید چالش برانگیز است؟ چرا؟ به چه طریقی می‌توانید دیتاست تولیدشده خود را چالش برانگیزتر و سخت‌تر کنید؟

ابتدا کتابخانه‌های مورد نیاز را import می‌کنیم.

```
1 from sklearn.datasets import load_breast_cancer, make_classification,
    make_blobs, make_circles, load_digits
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import pandas as pd
5 import plotly.express as px
6 from sklearn.decomposition import PCA
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LogisticRegression, SGDClassifier
9 from sklearn.preprocessing import StandardScaler, MinMaxScaler
10 from mpl_toolkits.mplot3d import Axes3D
11 from mlxtend.plotting import plot_decision_regions
```

Code 1: import libraries (Python)

سپس با استفاده از sklearn.datasets و مطابق ویژگی‌های مطرح شده توسط صورت سوال دیتاستی مناسب طبقه‌بندی و با ۱۰۰۰ نمونه و ۴ کلاس و ۳ ویژگی تولید می‌کنیم. برای تعداد نمونه‌ها با n-samples و برای تعداد کلاس‌ها با n-classes و برای تعداد ویژگی‌ها با n-features کار می‌کنیم و برابر با مقادیر خواسته شده در صورت سوال، قرارشان می‌دهیم. در این قسمت random-state را مطابق دورقم آخر شماره دانشجویی و برابر ۰۴ می‌گیریم. class-sep میزان تفکیک داده‌ها را از هم مشخص می‌کند. اگر مقدارش کم باشد، داده‌ها به هم نزدیک‌تر هستند و اگر بیشتر باشد، داده‌ها از هم دورتر می‌شوند. n-cluster-per-class تعداد خوشه‌های مربوط به هر کلاس را مشخص می‌کند. وقتی آن را برابر یک قرار می‌دهیم، یعنی هر کلاس یک خوشه جدا برای خود دارد.

```
1 # Generate dataset
2 X, y = make_classification(n_samples=1000,
3                             n_features=3,
```



```
4         n_classes=4,
5         n_redundant=0,
6         n_clusters_per_class=1,
7         class_sep=1 ,
8         random_state=4)
```

Code 2: Generate Dataset

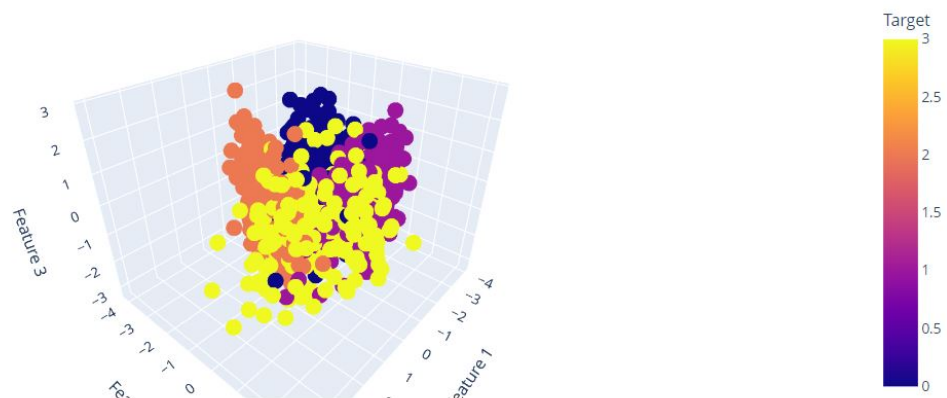
برای رسم مجموعه داده بصورت ۳D از کتابخانه `plotly.express` استفاده می کنیم. بدین منظور ابتدا داده های تولید شده را بصورت دیتافریم درمی آوریم و سه ویژگی و لیبل مخصوص هر داده را به عنوان ورودی به تابع مورد نظر می دهیم.

```
1 # Create a DataFrame from the generated dataset
2 df = pd.DataFrame(X, columns=['Feature 1', 'Feature 2', 'Feature 3'])
3 df['Target'] = y
4
5 # Plot the dataset using Plotly Express
6 fig = px.scatter_3d(df, x='Feature 1', y='Feature 2', z='Feature 3', color='
    Target', title='make_classification Dataset')
7 fig.show()
```

Code 3: Plot Dataset in 3D

خروجی مورد نظر بصورت زیر است. این نمودار بصورت سه بعدی بوده و امکان داشتن موقعیت مکانی هر داده و بررسی دقیق تر با بزرگنمایی را دارد.

make_classification Dataset



شکل ۲: 3D plot of dataset



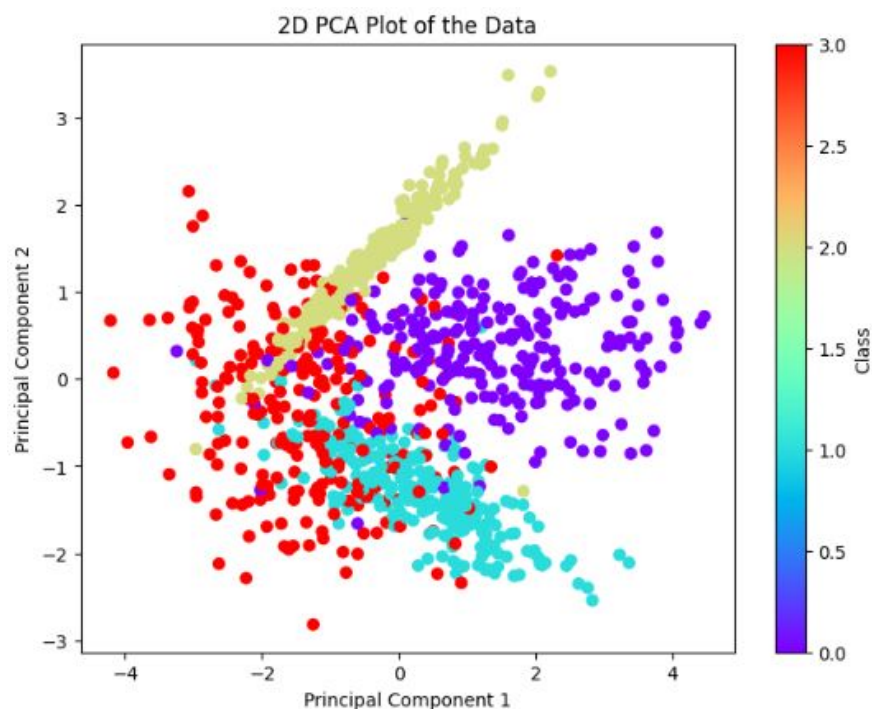
همچنین امکان رسم مجموعه داده بصورت دوبعدی و درک بهتر نیز فراهم است. برای این کار از کتابخانه `sklearn.decomposition` استفاده شده است. PCA می تواند کاهش بعد برای مجموعه داده را فراهمی کند بصورتی که با ادغام ویژگی های موجود در مجموعه داده و ترکیب آن ها با هم، از سه ویژگی حاضر در دیتاست آن را به دو ویژگی کاهش می دهد که قابل رسم در در فضای دوبعدی باشد.

```
1 # Perform PCA for dimensionality reduction
2 pca = PCA(n_components=2)
3 X_pca = pca.fit_transform(X)
4
5 # Plot the data
6 plt.figure(figsize=(8, 6))
7
8 # Scatter plot
9 plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap=plt.cm.rainbow)
10 plt.xlabel('Principal Component 1')
11 plt.ylabel('Principal Component 2')
12 plt.title('2D PCA Plot of the Data')
13
14 plt.colorbar(label='Class')
15
16 plt.show()
```

Code 4: Plot Dataset in 2D

خروجی این کد بصورت شکل ۳ زیر است.

چالش برانگیز بودن مجموعه داده تولید شده با توجه به تعداد کلاس ها و ویژگی های مشخص شده را می توان با مقادیر `n-clusters-per-class` که میزان چند قسمتی بودن هر کلاس را مشخص می کند و `class-sep` که میزان دوری یا درهم آمیختگی دیتاست را مشخص می کند، تنظیم کرد. در این مسئله این مقادیر برابر با ۱ قرار داده شده است. بنابراین از نظر چالش برانگیز بودن در حالت متوسط قرار دارد. اما می توان برای چالش برانگیز تر کردن `class-sep` را کاهش و یا `n-clusters-per-class` را افزایش داد. یک نمونه مثال در شکل ۴ آورده شده است. در این قسمت با افزایش تعداد `n-clusters-per-class` به مقدار ۲ و همچنین افزایش `class-sep` مقدار چالش برانگیز شدن دیتاست تغییر پیدا کرده است. اگر بخواهیم این چالش زیاد شود باشد مقدار این پارامترها را به ترتیب بزرگتر و کوچکتر کرد.



شکل ۳: 2D plot of dataset

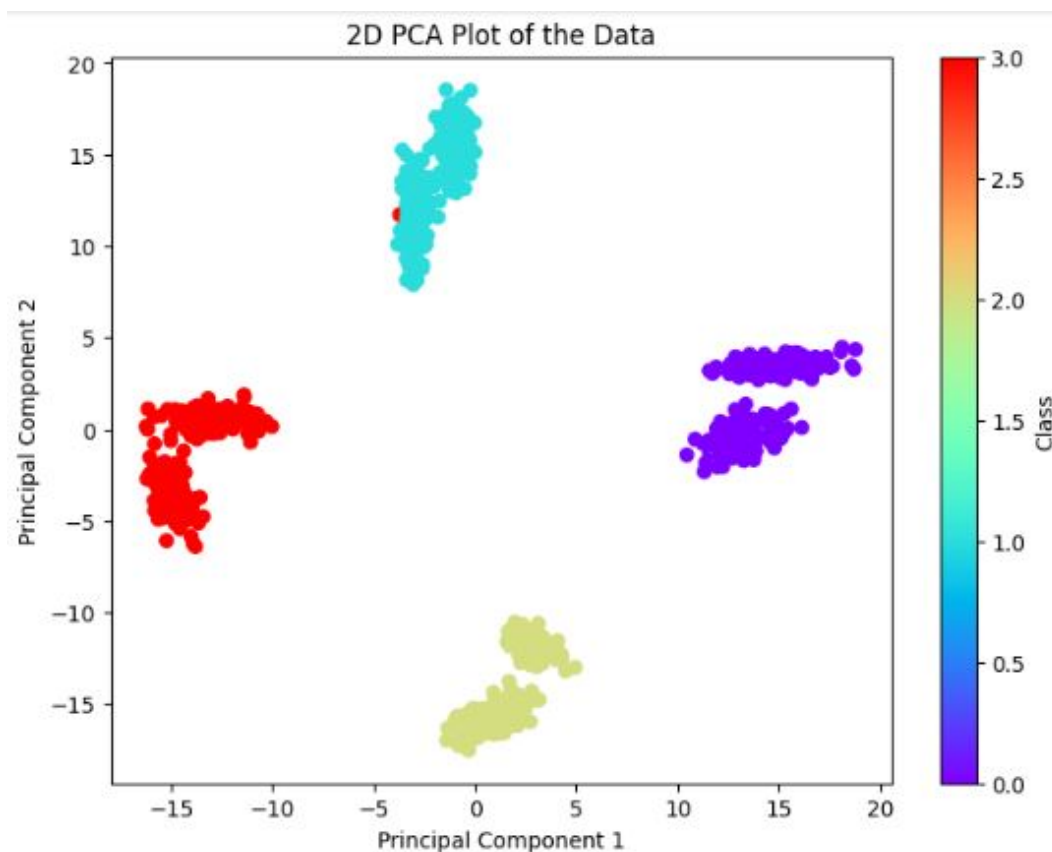
۳.۱ با استفاده از حداقل دو طبقه بند خطی آماده پایتون (در `model-linear.sklearn`) و در نظر گرفتن فرآیندهای مناسب، چهار کلاس موجود در دیتاست قسمت قبلی را از هم تفکیک کنید. ضمن توضیح روند انتخاب فرآیندها (مانند تعداد دوره آموزش و نرخ یادگیری)، نتیجه دقت آموزش و ارزیابی را نمایش دهید. برای بهبود نتیجه از چه تکنیک هایی استفاده کردید؟

در این مرحله ابتدا نیاز است که مجموعه داده خود را به دو قسمت آموزش و ارزیابی تقسیم کنیم. این مرحله بصورت آمده در زیر انجام می شود. مقدار ۲۰ درصد از داده ها را به قسمت ارزیابی اختصاص می دهیم. باید توجه داشت که این درصد بستگی به میزان داده های کل می تواند متفاوت باشد. هر چه تعداد کل داده ها بیشتر باشد، می شود درصد کمتری را به مجموعه داده های ارزیابی اختصاص داد. به عنوان مثال در یک مجموعه داده یک میلیونی میزان یک درصد از داده ها که برابر ۱۰ هزار داده است هم شاید برای مجموعه داده ارزیابی زیاد باشد و از این هم کمتر می شود اختصاص داد.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=4)
2 X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

Code 5: split dataset

در مرحله بعد لازم است که مقداری در مورد مجموعه داده خود اطلاعاتی بدست آوریم. بدین منظور با استفاده از کد زیر سعی شده است بازه تغییرات ویژگی های آموزش مورد بررسی قرار گیرد.



شکل ۴: n-clusters-per-class=2, class-sep=10

```
1 # Compute the range for each feature
2 feature_ranges = [(np.min(X_train[:, i]), np.max(X_train[:, i])) for i in
3                    range(X_train.shape[1])]
4
5 # Print the range for each feature
6 for i, (min_val, max_val) in enumerate(feature_ranges):
7     print(f"Feature {i+1}: Range = ({min_val}, {max_val})")
```

Code 6: feature range

خروجی نشان‌دهنده این است که بازه تغییرات هر سه ویژگی در یک محدوده مشخص است و تغییرات قابل توجه توجهی نسبت به یکدیگر ندارند.

به عنوان اولین طبقه‌بند آماده پایتون از LogisticRegression استفاده کرده می‌کنیم. این طبقه‌بند پارامترهای مختلفی را به عنوان ورودی می‌پذیرد. به عنوان مثال پارامتر solver نوع بهینه‌ساز را مشخص می‌کند و پارامتر C مقدار نرخ بهینه‌سازی را مشخص می‌کند. در این قسمت سعی می‌شود که به طور سازماندهی شده‌ای مقادیر مختلف برای این پارامترها انتخاب شود که بهینه‌ترین نتایج حاصل شود. بدین منظور از حلقه‌های تکرار برای بررسی مقادیر بهینه پارامترها استفاده می‌شود. پارامترهای در نظر گرفته شده برای تغییر در طول



```
Feature 1: Range = (-3.802236114367438, 3.716621174692759)
Feature 2: Range = (-4.126933452031901, 3.1082741888927683)
Feature 3: Range = (-3.2146554969049794, 2.8715878457267467)
```

شکل ۵: بازه تغییرات ویژگی‌های دیتاست

آموزش‌های پی‌درپی، C solver و tol هستند. تعداد دوره آموزش را هم با توجه به بررسی‌های صورت گرفته برابر ۱۰۰۰ قرار دادیم چون تعداد آموزش‌های بیشتر از آن به نتایج بهتر منجر نمی‌شد.

```
1 C_values = np.linspace(0.01, 1, 50)
2 solver_list = ['lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky', 'sag', '
   saga']
3 best_result = 0
4 for C in C_values:
5     for solver in solver_list:
6         for i in np.linspace(1e-5, 1e-3, 50):
7             LR_model = LogisticRegression(penalty='l2',
8                                             dual=False,
9                                             tol = i,
10                                            solver = solver,
11                                            max_iter=1000,
12                                            C = C,
13                                            random_state=4)
14             LR_model.fit(X_train, y_train)
15             result = LR_model.score(X_train, y_train)
16             if result > best_result:
17                 best_result = result
18                 print (f"best result is {best_result} with solver: {solver} , C= {C}
   and tol= {i}")
```

Code 7: train model

خروجی مورد نظر با مشخص بودن بهینه‌ترین پارامترها بصورت زیر است.
میزان نزدیکی داده‌ها به یکدیگر در مرحله تولید مجموعه داده باعث چالش برانگیز شدن مجموعه داده شده است. حال اگر مقدار class-sep را برابر ۲ بگذاریم که این به معنی فاصله بیشتر داده‌ها از یکدیگر است و سپس آموزش را انجام دهیم به نتایج زیر می‌رسیم.
این نتیجه نشان می‌دهد که اهمیت این پارامترها در تولید دیتاست و نتایج حاصل شده بالاست.
نتیجه آموزش و ارزیابی به صورت زیر قابل محاسبه است.



```
best result is 0.775 with solver: lbfgs , C= 0.01 and tol= 1e-05
best result is 0.78375 with solver: lbfgs , C= 0.030204081632653063 and tol= 1e-05
best result is 0.785 with solver: sag , C= 0.39387755102040817 and tol= 0.0005555102040816327
best result is 0.78625 with solver: sag , C= 0.5959183673469388 and tol= 0.00047469387755102045
best result is 0.7875 with solver: lbfgs , C= 0.9393877551020408 and tol= 1e-05
```

شکل ۶: نتیجه آموزش با پارامترهای مختلف

```
best result is 0.96875 with solver: lbfgs , C= 0.01 and tol= 1e-05
best result is 0.97125 with solver: lbfgs , C= 0.030204081632653063 and tol= 1e-05
best result is 0.9725 with solver: sag , C= 0.333265306122449 and tol= 0.0005555102040816327
```

شکل ۷: نتیجه آموزش با پارامترهای مختلف با $\text{class-sep}=2$

```
1 LR_model.fit(X_train, y_train)
2 LR_model.predict(X_test), y_test
3
4 print("LR_model train score = ", LR_model.score(X_train, y_train))
5 print("LR_model predict score = ", LR_model.score(X_test, y_test))
```

Code 8: train and test score

نتایج را به صورت زیر می توان دید.

```
LR_model train score = 0.9725
LR_model predict score = 0.98
```

شکل ۸: نتایج آموزش و ارزیابی

به عنوان دومین طبقه‌بند آماده پایتونی از SGDClassifier استفاده می‌کنیم. این طبقه‌بند نیز دارای پارامترهای ورودی مختلفی است. برای تست کردن اثر مقادیر مختلف پارامترهای و بررسی بهینه‌ترین آن‌ها مشابه قبل از حلقه‌های تکرار استفاده کرده‌ایم. ارامترهای در نظر گرفته شده برای تغییر در طول آموزش‌های پی‌درپی η که به معنی مقدار اولیه نرخ آموزشی است و نوع learning-rate که شامل دو نوع خودکار optimal و adaptive است و نوع تابع loss است. کد زیر به عنوان مرحله آموزش مورد استفاده قرار گرفته است.

```
1 loss_list = ['hinge', 'log_loss', 'squared_hinge', 'squared_error', 'huber']
2 learning_rate = ['optimal', 'adaptive']
3 best_result = 0
4 for loss in loss_list:
5     for lr in learning_rate:
```




```

6 for eta in np.linspace(1e-3, 30, 50):
7     SGD_model = SGDClassifier(loss=loss,
8                               max_iter=2000,
9                               tol=1e-3,
10                              eta0=eta,
11                              learning_rate=lr,
12                              random_state=4)
13     SGD_model.fit(X_train, y_train)
14     result = SGD_model.score(X_train, y_train)
15     if result > best_result:
16         best_result = result
17     print(f"best result is {best_result} with loss: {loss} , learning
rate= {lr} and eta0= {eta}")

```

Code 9: SGDClassifier

نتایج آموزش و با خروجی گرفتن بهینه‌ترین مقادیر پارامترهای موجود به صورت زیر است.

```

best result is 0.95875 with loss: hinge , learning rate= optimal and eta0= 0.001
best result is 0.9625 with loss: hinge , learning rate= adaptive and eta0= 0.001
best result is 0.96625 with loss: hinge , learning rate= adaptive and eta0= 0.6132244897959184
best result is 0.9675 with loss: hinge , learning rate= adaptive and eta0= 6.735469387755102

```

شکل ۹: نتایج آموزش با پارامترهای مختلف

کد مورد استفاده برای نتایج آموزش و ارزیابی و خروجی مورد نظر بصورت زیر آمده است.

```

1 SGD_model.fit(X_train, y_train)
2 SGD_model.predict(X_test), y_test
3
4 print("SGD_model train score = ", SGD_model.score(X_train, y_train))
5 print("SGD_model predict score = ", SGD_model.score(X_test, y_test))

```

Code 10: SGD train and test score

```

SGD_model train score = 0.9675
SGD_model predict score = 0.975

```

شکل ۱۰: نتایج آموزش و تست SGD



برای بهبود نتیجه در مرحله اول سعی شد که با استفاده از تغییرات پارامترهای مورد استفاده در هر طبقه‌بند، بهینه‌ترین نتایج حاصل شود. بدین منظور ابتدا تعداد دوره‌های آموزشی افزایش داده شد که به خوبی داده‌های مختلف در شبکه آموزش ببینند و سپس با استفاده از حلقه‌های تکرار سعی شد بهینه‌ترین مقدار برای سایر پارامترها انتخاب شود. مورد دیگری که مورد بررسی قرار گرفت میزان تغییرات ویژگی‌های دیتاست بود. همانطور که پیش از این بررسی شد، تغییرات در مورد این دیتاست به خصوص، دارای تفاوت زیادی در بین ویژگی‌های خود نبود و این یک ویژگی مناسب برای یک دیتاست است. اگر دیتاستی دارای ویژگی‌های با بازه تغییرات بسیار متفاوت باشد، نیاز است ابتدا نرمال‌سازی داده صورت بگیرد. البته در این مرحله نرمال‌سازی داده صورت گرفت و در کد موجود است اما به دلیل اینکه بازه تغییرات خود مجموعه داده مناسب بود و با اعمال نرمال‌سازی با تابع `minmaxscaler` نتایج بهتری حاصل نشد، از آوردن در گزارش کار خودداری شد.

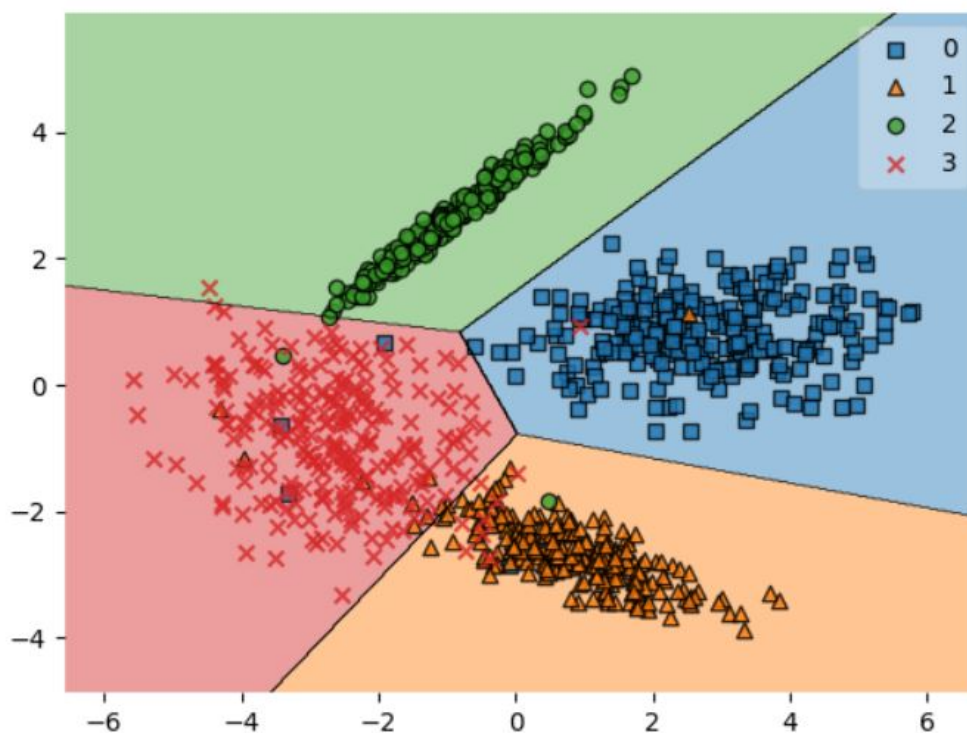


۴.۱ مرز و نواحی تصمیم گیری برآمده از مدل آموزش دیده خود را به همراه نمونه ها در یک نمودار نشان دهید. اگر می توانید نمونه هایی که اشتباه طبقه بندی شده اند را با شکل و رنگ متفاوت نمایش دهید.

در این قسمت برای رسم داده ها به همراه مرز تصمیم گیری هر طبقه از کد آماده پایتونی کتابخانه `mlxtend` استفاده شده است. کد مورد استفاده و نمودار مورد نظر بصورت زیر هستند. همانطور که قابل مشاهده است داده های به طور مناسبی داخل مرزهای تصمیم گیری قرار دارند.

```
1 from mlxtend.plotting import plot_decision_regions
2 plot_decision_regions(X_pca, y, clf=LR_model_pca)
```

Code 11: Import drawdata library



شکل ۱۱: نتایج آموزش و تست SGD



۵.۱ فرآیندی مشابه قسمت «۲» را با تعداد کلاس و ویژگی دلخواه؛ اما با استفاده از ابزار drawdata تکرار کنید. قسمت های «۳» و «۴» را برای این داده های جدید تکرار و نتایج را به صورتی مناسب نشان دهید. ابتدا لازم است کتابخانه مورد نیاز برای استفاده از این ابزار را فراخوانی کنیم.

```
1 !pip install drawdata
2 from drawdata import ScatterWidget
```

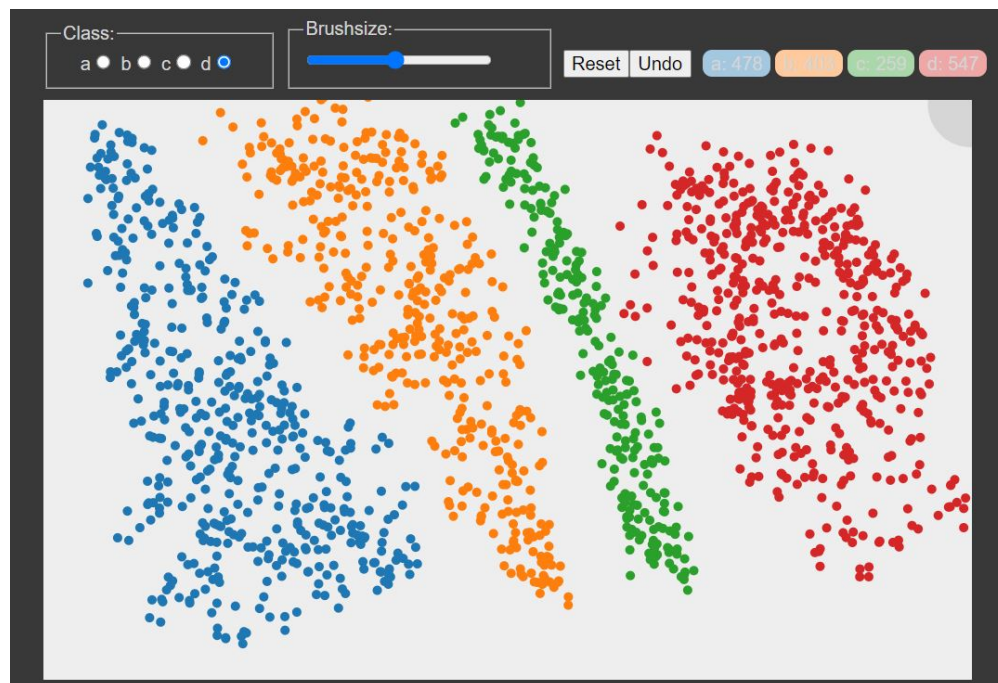
Code 12: Import drawdata library

با استفاده از کد زیر می توان ابزار مورد نظر را فراخوانی کرد.

```
1 widget = ScatterWidget()
2 widget
```

Code 13: drawdata

پس از فراخوانی با ایجاد پنجره ای مشابه زیر دیتاست مورد نظر را طراحی می کنیم.



شکل ۱۲: دیتاست طراحی شده با drawdata

در این قسمت ابتدا دیتاست طراحی شده را بصورت یک دیتافریم درآورده و مقادیر عددی برای برچسب های این داده را متناسب با طبقه آن دیتا به آن اختصاص می دهیم. همچنین دو ستون ابتدایی این مجموعه داده را که ویژگی های آن هستند را به مقادیر X و ستون انتهایی را به Y عنوان برچسب اختصاص می دهیم.



```
1 # Get the drawn data as a list of dictionaries
2 widget.data
3 # Get the drawn data as a dataframe
4 draedata_df = widget.data_as_pandas
5
6 draedata_df.replace({'a': 1, 'b': 2, 'c': 3, 'd': 4}, inplace=True)
7 draedata_df
8
9 X = draedata_df.iloc[:, :2].values
10 y = draedata_df.iloc[:, 3].values
```

Code 14: creat dataframe

نیاز است که پیش از آموزش مجموعه داده موجود را به دو بخش آموزش و ارزیابی تقسیم کنیم. این کار بصورت زیر انجام می پذیرد. با استفاده از کد موجود در خط دوم نیز می توان به تعداد هر کدام از مجموعه داده های تقسیم شده پی برد.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=4)
2 X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

Code 15: Split Dataset

به عنوان اولین طبقه بند از LogisticRegression استفاده می کنیم. در این قسمت از مقادیر بهینه ای که کتابخانه sklearn بصورت پیش فرض پیشنهاد داده است استفاده شده است.

```
1 LR_model = LogisticRegression(penalty='l2',
2                                dual=False,
3                                tol = 1e-05,
4                                solver='lbfgs',
5                                max_iter=1000,
6                                C=1,
7                                random_state=4)
8 LR_model.fit(X_train, y_train)
9 LR_model.predict(X_test)
10 print("LR_model train score = ", LR_model.score(X_train, y_train))
11 print("LR_model test score = ", LR_model.score(X_test, y_test))
```

Code 16: Train by LR



نتیجه آموزش و ارزیابی بصورت زیر قابل مشاهده است.

```
LR_model train score = 1.0
LR_model test score = 1.0
```

شکل ۱۳: نتیجه آموزش و ارزیابی با استفاده از Regression Logistic

برای آموزش SGD از حلقه‌های تکرار برای پیدا کردن بهینه‌ترین پارامترها استفاده کردیم.

```
1 loss_list = ['hinge', 'log_loss', 'squared_hinge', 'squared_error', 'huber']
2 learning_rate = ['optimal', 'adaptive']
3 best_result = 0
4 for loss in loss_list:
5     for lr in learning_rate:
6         for eta in np.linspace(1e-3, 30, 50):
7             SGD_model = SGDClassifier(loss=loss,
8                                     max_iter=2000,
9                                     tol=1e-3,
10                                    eta0=eta,
11                                    learning_rate=lr,
12                                    random_state=4)
13             SGD_model.fit(X_train, y_train)
14             result = SGD_model.score(X_train, y_train)
15             if result > best_result:
16                 best_result = result
17                 print(f"best result is {best_result} with loss: {loss} , learning
rate= {lr} and eta0= {eta}")
```

Code 17: Train by SGD

و نتایج حاصله به صورت زیر درآمد.

```
best result is 0.396590066716086 with loss: hinge , learning rate= optimal and eta0= 0.001
best result is 0.5982209043736101 with loss: hinge , learning rate= adaptive and eta0= 0.001
best result is 0.8495181616011861 with loss: hinge , learning rate= adaptive and eta0= 0.6132244897959184
best result is 0.8628613787991104 with loss: hinge , learning rate= adaptive and eta0= 1.2254489795918366
best result is 0.8747220163083765 with loss: hinge , learning rate= adaptive and eta0= 3.0621224489795917
best result is 0.8776871756856931 with loss: hinge , learning rate= adaptive and eta0= 5.511020408163266
best result is 0.8888065233506302 with loss: hinge , learning rate= adaptive and eta0= 14.69438775510204
best result is 0.8977020014825797 with loss: squared_hinge , learning rate= adaptive and eta0= 0.001
```

شکل ۱۴: نتیجه آموزش با استفاده از SGD



در نهایت نتایج آموزش و ارزیابی به صورت زیر قابل مشاهده است.

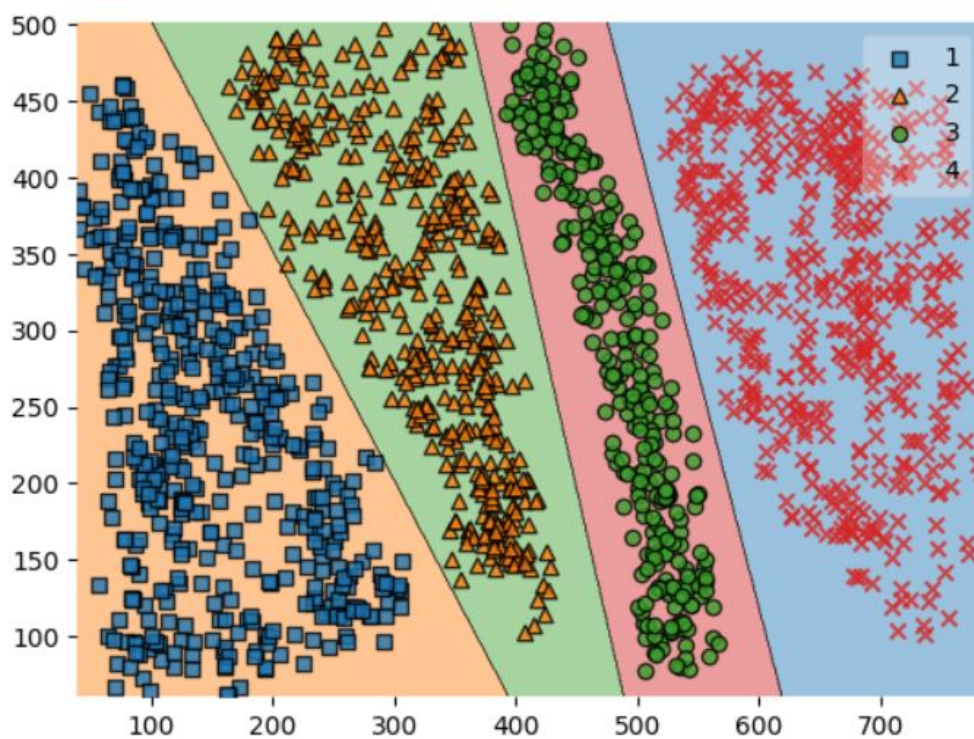
```
SGD_model train score = 0.8977020014825797  
SGD_model test score = 0.8816568047337278
```

شکل ۱۵: نتیجه آموزش و ارزیابی با استفاده از SGD

نمودار نهایی داده‌ها به همراه مرز تصمیم‌گیری در این قسمت بصورت استفاده از کتابخانه آماده پایتونی مشابه قبل صورت می‌گیرد.

```
1 from mlxtend.plotting import plot_decision_regions  
2 plot_decision_regions(X, y, clf=LR_model)
```

Code 18: Train by SGD



شکل ۱۶: نتیجه آموزش و ارزیابی با استفاده از SGD



۲ عنوان سوال دوم

۱.۲ با مراجعه به صفحه دیتاست CWRU Bearing با یک دیتاست مربوط به حوزه «تشخیص عیب» آشنا شوید. با جستجوی آن در اینترنت و مقالات، توضیحاتی از اهداف، ویژگی ها و حالت های مختلف این دیتاست ارائه کنید.

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 ۰.۰۷ تا ۰.۴۰ 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 خطا)، 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 ۰.۰۷ 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]

۲.۲ قسمت دوم

۱.۲.۲ از هر کلاس نمونه با طول N جدا کنید (حداقل ۱۰۰ و N حداقل ۲۰۰ باشد). یک ماتریس از داده های هر دو کلاس به همراه برجسب مربوطه تشکیل دهید. می توانید پنجره ای به طول N در نظر بگیرید و در نهایت یک ماتریس $N \times$ از داده های هر کلاس استخراج کنید.

ابتدا مجموعه داده های مورد خواسته مسئله و کتابخانه های مورد نیاز را فراخوانی می کنیم. سپس دو 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 دستاها دسترسی پیدا کرد.

```
1 !gdown 14JXJeFLn-V1KSM5yNAiwokaX1F9P3haS
```

```
2 !gdown 1oqLyBpuEfX0pYIZ-aY_gzznCPsBz1pGp
```

```
3
```

```
4 from scipy.io import loadmat
```

```
5 import numpy as np
```




```

6 from scipy.stats import skew, kurtosis
7 from sklearn.utils import shuffle
8 from sklearn.model_selection import train_test_split
9
10 # Load the .mat file
11 normal_data = loadmat('/content/97.mat')
12 fault_data = loadmat('/content/105.mat')
13
14 print(normal_data.keys())
15 print(fault_data.keys())

```

Code 19: import dataset and library

خروجی سرستون‌های دیتاست‌ها بصورت زیر 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'])

```

شکل ۱۷: نتیجه آموزش و ارزیابی با استفاده از SGD

سپس از مجموعه داده سالم ستون X097-DE-time و از مجموعه داده خطا ستون X105-DE-time را به عنوان نمونه با توجه به خواسته مسئله انتخاب می‌کنیم و در متغیرهای مناسب می‌ریزیم. برای داشتن دید مناسب از ابعاد مجموعه داده ابتدا نوع ستون‌های مورد نظر و همچنین ابعاد آن‌ها را خروجی می‌گیریم.

```

1 normal_data_variable = normal_data['X097_DE_time']
2 print(type(normal_data_variable)) # Print the type of the variable
3 print(normal_data_variable.shape) # Print the shape of the variable (if it's
  a numpy array)
4
5 fault_data_variable = fault_data['X105_DE_time']
6 print(type(fault_data_variable)) # Print the type of the variable
7 print(fault_data_variable.shape) # Print the shape of the variable (if it's a
  numpy array)

```

Code 20: import dataset and library

خروجی مورد نظر به صورت زیر است.

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) می‌شود.



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

شکل ۱۸: نتیجه آموزش و ارزیابی با استفاده از SGD

```
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 = fault_data_variable[:n_samples * n_data, 0]
6
7 print(ext_normal_data.shape)
8 print(ext_fault_data.shape)
```

Code 21: import dataset and library

سپس برا تکمیل ابعاد مجموعه داده با استفاده از دستور reshape، دیتاها به صورتی درمی آیند که سطرهاى آن دربرگیرنده ۱۲۰ نمونه مورد خواسته مسئله باشد.

```
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 = ext_fault_data.reshape(n_samples, n_data)
4
5 print(ext_normal_data.shape)
6 print(ext_fault_data.shape)
```

Code 22: import dataset and library

خروجی این کد نیز بصورت زیر است.

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

شکل ۱۹: نتیجه آموزش و ارزیابی با استفاده از SGD



۲.۲.۲ در مورد اهمیت استخراج ویژگی در یادگیری ماشین توضیحاتی بنویسید. سپس، با استفاده از حداقل ۸ عدد از روش های ذکر شده در جدول ۱، ویژگی های دیتاست قسمت «۲ آ» را استخراج کنید و یک دیتاست جدید تشکیل دهید.

استخراج ویژگی فرآیند شناسایی و انتخاب مهم ترین اطلاعات یا ویژگی ها از یک مجموعه داده است. این مانند تقطیر عناصر ضروری است، کمک به ساده سازی و برجسته کردن جنبه های کلیدی و در عین حال فیلتر کردن جزئیات کمتر مهم. این راهی برای تمرکز بر آنچه واقعاً در داده ها مهم است.

استخراج ویژگی مهم است زیرا اطلاعات پیچیده را ساده تر می کند. در مواردی مانند یادگیری کامپیوتری، به یافتن مهم ترین الگوها یا جزئیات کمک می کند و با تمرکز بر آنچه در داده ها اهمیت دارد، رایانه ها را در پیش بینی یا تصمیم گیری بهتر می کند.

اهمیت استخراج ویژگی در یادگیری ماشین را می توان بصورت خلاصه اینطور بررسی کرد:

کاهش ابعاد: این کار باعث ساده تر شدن داده های با ابعاد بالا می شود، که محاسبات سریع تری را فراهم می کند و از بیش برآزش (overfitting) جلوگیری می کند.

کاهش نویز: این کار با حذف ویژگی های بی اهمیت یا نویزی، عملکرد مدل را بهبود می بخشد.

بهبود عملکرد: با تمرکز بر ویژگی های مرتبط، مدل ها می توانند الگوهای داده را بهتری تشخیص دهند که منجر به دقت و عمومیت بهتر می شود.

قابل فهم بودن: استخراج ویژگی می تواند مدل ها را قابل فهم تر کند با تبدیل داده به یک فرمت قابل فهم تر.

کار با داده های غیر عددی: این اجازه می دهد از انواع داده های غیر عددی (مانند متن یا تصاویر) با تبدیل آن ها به یک فرمت مناسب برای الگوریتم های یادگیری ماشین استفاده شود.

ادغام دانش حوزه ای: استخراج ویژگی اغلب شامل دانش خاص حوزه است که می تواند اطلاعاتی را تقویت کند.

با توجه به جدول مسئله تعداد هشت عدد از روش های ذکر شده انتخاب شد. با استفاده از این این روش ها ویژگی های دیتاست مورد استخراج می شود. کد زیر نحوه محاسبه هر کدام از این روش ها و اعمال بر روی دیتاست نرمال را نشان می دهد. برای محاسبه هر کدام از این روش ها از کتابخانه numpy استفاده شده است. قسمت axis=1 باعث می شود که الگوریتم مورد نظر بر روی سطرها که نشان دهنده نمونه های دیتاست هستند اعمال شود. خروجی ابعاد هر کدام از الگوریتم در پایان پرینت شده است. در قسمت بعدی خروجی ویژگی های استخراج شده با دستور column-stack بصورت ستونی به یکدیگر متصل شدند تا دیتاستی با ابعاد (۸، ۱۲۰) ساخته شود. دی قسمت پایانی نیز یک ستون صفر به دیتاست اضافه شده است که نشان دهنده برچسب کلاس نرمال است.

```
1 #Normal data feature extraction
2 normal_standard_deviations = np.std(ext_normal_data, axis=1)
3 normal_skewnesses = skew(ext_normal_data, axis=1)
4 normal_kurtoses = kurtosis(ext_normal_data, axis=1)
5 normal_peak_to_peaks = np.ptp(ext_normal_data, axis=1)
6 normal_root_mean_squares = np.sqrt(np.mean(np.square(ext_normal_data), axis=1)
7 )
8 normal_means = np.mean(ext_normal_data, axis=1)
9 normal_absolute_means = np.mean(np.abs(ext_normal_data), axis=1)
10 normal_peaks = np.max(ext_normal_data, axis=1)
11 print("Standard Deviations:", normal_standard_deviations.shape)
```



```

12 print("Skewnesses:", normal_skewnesses.shape)
13 print("Kurtoses:", normal_kurtoses.shape)
14 print("Peak to Peaks:", normal_peak_to_peaks.shape)
15 print("Root Mean Squares:", normal_root_mean_squares.shape)
16 print("Means:", normal_means.shape)
17 print("Absolute Means:", normal_absolute_means.shape)
18 print("Peaks:", normal_peaks.shape)
19
20 normal_ext_feature_dataset = np.column_stack((normal_standard_deviations,
        normal_skewnesses, normal_kurtoses, normal_peak_to_peaks,
21        normal_root_mean_squares, normal_means,
        normal_absolute_means, normal_peaks))
22
23
24 normal_ext_feature_dataset = np.hstack((normal_ext_feature_dataset, np.zeros((
        len(normal_ext_feature_dataset), 1))))
25 print("final normal dataset: ", normal_ext_feature_dataset.shape)

```

Code 23: import dataset and library

خروجی کد مورد استفاده بصورت زیر است. در نهایت دیتاست مربوط به داده‌های نرمال با ابعاد (۹، ۱۲۰) که هم دارای ویژگی‌ها است و هم دارای برچسب در ستون آخر خود، قابل بهره‌برداری است.

```

Standard Deviations: (120,)
Skewnesses: (120,)
Kurtoses: (120,)
Peak to Peaks: (120,)
Root Mean Squares: (120,)
Means: (120,)
Absolute Means: (120,)
Peaks: (120,)
final normal dataset: (120, 9)

```

شکل ۲۰: نتیجه آموزش و ارزیابی با استفاده از SGD

برای دیتاست خطا نیز همین مسیر مورد استفاده قرار گرفته است که کد و خروجی آن بصورت زیر قابل مشاهده است.

```

1 #fault data feature extraction
2 fault_standard_deviations = np.std(ext_fault_data, axis=1)
3 fault_skewnesses = skew(ext_fault_data, axis=1)
4 fault_kurtoses = kurtosis(ext_fault_data, axis=1)
5 fault_peak_to_peaks = np.ptp(ext_fault_data, axis=1)

```



```
6 fault_root_mean_squares = np.sqrt(np.mean(np.square(ext_fault_data), axis=1))
7 fault_means = np.mean(ext_fault_data, axis=1)
8 fault_absolute_means = np.mean(np.abs(ext_fault_data), axis=1)
9 fault_peaks = np.max(ext_fault_data, axis=1)
10
11 print("Standard Deviations:", fault_standard_deviations.shape)
12 print("Skewnesses:", fault_skewnesses.shape)
13 print("Kurtoses:", fault_kurtoses.shape)
14 print("Peak to Peaks:", fault_peak_to_peaks.shape)
15 print("Root Mean Squares:", fault_root_mean_squares.shape)
16 print("Means:", fault_means.shape)
17 print("Absolute Means:", fault_absolute_means.shape)
18 print("Peaks:", fault_peaks.shape)
19
20 fault_ext_feature_dataset = np.column_stack((fault_standard_deviations,
        fault_skewnesses, fault_kurtoses, fault_peak_to_peaks,
21        fault_root_mean_squares, fault_means,
        fault_absolute_means, fault_peaks))
22
23 fault_ext_feature_dataset = np.hstack((fault_ext_feature_dataset, np.ones((len
        (fault_ext_feature_dataset), 1))))
24 print("final fault dataset: ", fault_ext_feature_dataset.shape)
```

Code 24: import dataset and library

```
Standard Deviations: (120,)
Skewnesses: (120,)
Kurtoses: (120,)
Peak to Peaks: (120,)
Root Mean Squares: (120,)
Means: (120,)
Absolute Means: (120,)
Peaks: (120,)
final fault dataset: (120, 9)
```

شکل ۲۱: نتیجه آموزش و ارزیابی با استفاده از SGD

در نهایت برای این که دیتاست تشکیل شده مانند یک دیتاستی که دارای کلاس‌های مختلف است به صورت استاندارد دربیاید، دو 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,
    fault_ext_feature_dataset), axis=0)
2 final_data.shape

```

Code 25: import dataset and library

(240, 9)

شکل ۲۲: نتیجه آموزش و ارزیابی با استفاده از SGD

۳.۲.۲ ضمن توضیح اهمیت فرآیند ب رزدن (مخلوط کردن) داده ها را در صورت امکان مخلوط کرده و با نسبت تقسیم دلخواه و معقول به دو بخش «آموزش» و «ارزیابی» تقسیم کنید.

مخلوط کردن داده ها می تواند منجر به بهبود دقت مدل های پیش بینی و یادگیری ماشین شود. این فرآیند اغلب منجر به داده هایی با تنوع و اطلاعات بیشتری می شود که باعث افزایش دقت و قابلیت پیش بینی مدل های آموزش داده شده می شود. ادغام داده ها اغلب منجر به افزایش تنوع و کمیت داده های موجود می شود که این امر می تواند منجر به بهتر شدن تحلیل ها، کاهش انحرافات غیرمنتظره و افزایش اطمینان در تصمیم گیری ها شود. مخلوط کردن داده ها ممکن است اطلاعات جدید و مفیدی که از یک ترکیب منابع به دست می آید را ارائه دهد که ممکن است در انجام تحلیل های پیشرفته و تصمیم گیری های مهم مورد استفاده قرار گیرد.

بنابراین مخلوط کردن دیتاها دارای اهمیت فراوانی است. در این قسمت پس از مخلوط کردن دیتاها با دستور shuffle، دیتاهای مربوط به ویژگی های دیتاست (همه ستون ها به جز ستون آخر) به متغیر X و دیتای برچسب به متغیر y تعلق پیدا می کند. در نهایت دیتاها به دو دسته آموزش و ارزیابی تقسیم می شوند. نسبت داده های آموزش به ارزیابی در این قسمت با توجه به اندازه کلی دیتاست برابر ۲۰ درصد در نظر گرفته شده است. 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
6 x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.1)
7 x_train.shape, x_test.shape, y_train.shape, y_test.shape

```

Code 26: import dataset and library

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.



```
((216, 8), (24, 8), (216, 1), (24, 1))
```

شکل ۲۳: نتیجه آموزش و ارزیابی با استفاده از SGD

```
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 27: import dataset and library

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

```
((216, 8), (24, 8), (216, 1), (24, 1))
```

شکل ۲۴: نتیجه آموزش و ارزیابی با استفاده از SGD

۴.۲.۲ حداقل دو روش برای نرمال سازی داده ها را با ذکر اهمیت این فرآیند توضیح دهید و با استفاده از یکی از این روش ها، داده ها را نرمال کنید. آیا از اطلاعات بخش «ارزیابی» در فرآیند نرمال سازی استفاده کردید؟ چرا؟

نرمال کردن داده ها برای منطقی کردن مقایسه پذیری و بهبود عملکرد الگوریتم های یادگیری ماشین و تسهیل فرآیند بهینه سازی و کمک به جلوگیری از overfitting صورت می گیرد. واحدها و مقیاس های مختلف برای ویژگی ها ممکن است باعث کاهش دقت و عدم قابلیت مقایسه و تفسیر شود. نرمال کردن داده ها باعث می شود تمام ویژگی ها به یک مقیاس یا بازه مشابه تبدیل شوند که قابل مقایسه تر و تفسیر پذیرتر باشند. برای مدل های با پارامتر زیاد، اگر داده ها نرمال نشوند، احتمال overfitting بیشتر می شود؛ زیرا مدل ممکن است به اندازه ی زیادی به داده های با ویژگی های بزرگتر (به دلیل مقیاس بزرگتر) وابستگی پیدا کند و از یادگیری الگوهای کلی دور شود. روش اول Scaling Max-Min: در این روش، داده ها به گونه ای تغییر می کنند که حداقل و حداکثر آنها به ترتیب به یک مقدار نگاشته می شوند. فرمول نرمال سازی Max-Min برای یک داده X به صورت زیر است:

$$X_{\text{normalized}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \quad (۱)$$

روش دوم Normalization score-Z: این روش با توجه به میانگین و انحراف معیاری داده ها، آن ها را به صورتی نرمال می کند که میانگین آن ها صفر و انحراف معیاری یک باشد.

$$X_{\text{normalized}} = \frac{X - \mu}{\sigma} \quad (۲)$$

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 فرآیند نرمال سازی با استفاده از اطلاعات داده ها آموزش بر روی داده های آموزش و ارزیابی صورت می پذیرد.



برای نرمال‌سازی داده‌ها در این سوال از روش اول استفاده می‌کنیم. کد مورد استفاده بصورت زیر است.

```
1 from sklearn.preprocessing import MinMaxScaler, StandardScaler
2
3 # Initialize MinMaxScaler
4 scaler = MinMaxScaler()
5
6 # Fit the scaler on the training data
7 scaler.fit(x_train)
8
9 # Transform the training and testing data
10 x_train = scaler.transform(x_train)
11 x_test = scaler.transform(x_test)
```

Code 28: import dataset and library

۳.۲ بدون استفاده از کتابخانه‌های آماده پایتون، مدل طبقه‌بند، تابع اتلاف و الگوریتم یادگیری و ارزیابی را کدنویسی کنید تا دو کلاس موجود در دیتاست به خوبی از یکدیگر تفکیک شوند. نمودار تابع اتلاف را رسم کنید و نتیجه ارزیابی روی داده‌های تست را با حداقل ۲ شاخصه محاسبه کنید. نمودار تابع اتلاف را تحلیل کنید. آیا می‌توان از روی نمودار تابع اتلاف و قبل از مرحله ارزیابی با قطعیت در مورد عمل کرد مدل نظر داد؟ چرا و اگر نمی‌توان، راه حل چیست؟

طبق تابع‌های موجود برای ایجاد \hat{y} با استفاده از تابع sigmoid به شکل زیر عمل می‌کنیم و ابتدا تابع sigmoid را تعریف می‌کنیم و سپس \hat{y} را ایجاد می‌کنیم. با استفاده از \hat{y} تابع logistic regression که ضرب x ها و w هاست را تشکیل می‌دهیم و خروجی آن \hat{y} می‌شود که همان y است که ما در شبکه خود ایجاد کرده ایم. سپس تابع اتلاف را با $\text{loss} = \log$ معرفی کرده و گرادین آن را محاسبه می‌کنیم و η که همان ضریب یادگیری است را در grads که حاصل ضرب x در اختلاف y اصلی ما با y بدست آمده است، ضرب کرده و w را هر سری با آن آپدیت می‌کنیم تا به w درست برسیم

```
1 def sigmoid(x):
2     return 1 / (1 + np.exp(-x))
3
4 def logistic_regression(x, w):
5     y_hat = sigmoid(x @ w)
6     return y_hat
```




```

7
8 def bce(y, y_hat):
9     loss = -(np.mean(y*np.log(y_hat) + (1-y)*np.log(1-y_hat)))
10    return loss
11
12 def gradient(x, y, y_hat):
13     grads = (x.T @ (y_hat - y)) / len(y)
14     return grads
15
16 def gradient_descent(w, eta, grads):
17     w -= eta*grads
18     return w
19
20 def accuracy(y, y_hat):
21     acc = np.sum(y == np.round(y_hat))/len(y)
22     return acc

```

Code 29: import dataset and library

در ادامه برای در نظر گرفتن بایاس یک ستون تماماً یک به ان اضافه می‌کنیم. در اینجا خروجی گرفتن ابعاد دیتاست دارای اهمیت است چون با هر بار اجرای این سلول سک ستون به دیتا اضافه می‌شود. کد و خروجی مورد نظر بصورت زیر است.

```

1 x_train = np.hstack((np.ones((len(x_train), 1)), x_train))
2 x_train.shape

```

Code 30: import dataset and library

(216, 9)

شکل ۲۵: نتیجه آموزش و ارزیابی با استفاده از SGD

در این قسمت نیاز است تا پیش از شروع آموزش برخی پارامترهای مورد نیاز تعریف شوند. مقداردهی اولیه وزن‌ها با توجه به تعداد ویژگی‌ها دیتاست و با در نظر گرفتن بایاس صورت می‌گیرد. پارامترهای مورد نیاز تعریف و یک لیست خالی برای ذخیره کردن مقدار خطا در هر مرحله ایجاد می‌شود. سپس در یک حلقه تکرار فرآیند آموزش اجرا می‌شود.

```

1 m = 8 #num of features

```



```

2 w = np.random.randn(m+1, 1)
3
4 num_epoch = 1000
5 eta = 0.01
6
7 error_hist = []
8
9 for epoch in range(num_epoch):
10     y_hat = logistic_regression(x_train, w)
11     loss = bce(y_train, y_hat)
12     error_hist.append(loss)
13     grads = gradient(x_train, y_train, y_hat)
14     w = gradient_descent(w, eta, grads)
15     if epoch % 50 == 0:
16         print(f'Epoch = {epoch+1}, \t loss = {loss:.4f}, \t w={w.T[0, 0]}')

```

Code 31: import dataset and library

خروجی آموزش به صورت زیر گزارش شده است. کد طوری نوشته شده است که هر ۵۰ اپیاک یکبار خروجی تعیین شده نمایان شود. همانطور که قابل مشاهده است مقدار خطا در هر مرحله دارای شیب کاهشی است که نشان دهنده فرآیند مناسب آموزش است. وزن‌ها نیز در هر مرحله دارای تغییر و حرکت به سمت مقادیر بهینه‌اند.

```

Epoch = 1,      loss = 1.705,    w=-0.8901391594909194
Epoch = 51,     loss = 1.268,    w=-0.7314287370227337
Epoch = 101,    loss = 0.9389,   w=-0.6116008682208355
Epoch = 151,    loss = 0.7309,   w=-0.5384661548822829
Epoch = 201,    loss = 0.6119,   w=-0.5077540513480393
Epoch = 251,    loss = 0.5418,   w=-0.5078634612413715
Epoch = 301,    loss = 0.4957,   w=-0.5279823807819595
Epoch = 351,    loss = 0.4616,   w=-0.5605434166192694
Epoch = 401,    loss = 0.4342,   w=-0.6006668044127286
Epoch = 451,    loss = 0.4107,   w=-0.6452400506262735
Epoch = 501,    loss = 0.39,     w=-0.6922583629886134
Epoch = 551,    loss = 0.3714,   w=-0.7404144533029205
Epoch = 601,    loss = 0.3544,   w=-0.7888486819568431
Epoch = 651,    loss = 0.3387,   w=-0.8369946942048667
Epoch = 701,    loss = 0.3242,   w=-0.884481770148694
Epoch = 751,    loss = 0.3108,   w=-0.9310715404646281
Epoch = 801,    loss = 0.2983,   w=-0.9766160796975121
Epoch = 851,    loss = 0.2867,   w=-1.0210296489878985
Epoch = 901,    loss = 0.2758,   w=-1.0642693611264002
Epoch = 951,    loss = 0.2657,   w=-1.106321796993851

```

شکل ۲۶: نتیجه آموزش و ارزیابی با استفاده از SGD

برای رسم نمودار تابع اتلاف از دستور plt=plot استفاده می‌کنیم و نتیجه آموزش نیز با استفاده از تابع تعریف شده در قسمت قبل



دارای محاسبه است.

```
1 import matplotlib.pyplot as plt
2 acc_train = accuracy(y_train, y_hat)
3 plt.plot(error_hist)
```

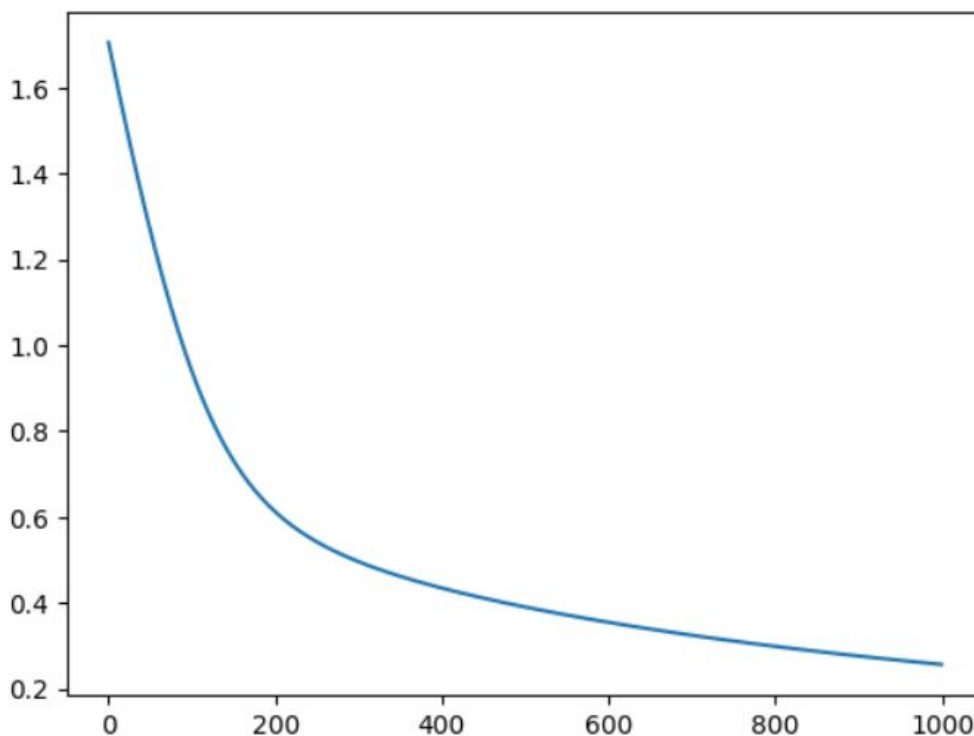
Code 32: import dataset and library

مقدار acc برای داده‌های آموزش بصورت زیر است.

```
train accuracy = 1.0
```

شکل ۲۷: نتیجه آموزش و ارزیابی با استفاده از SGD

همچنین نمودار تابع اتلاف بصورت زیر قابل مشاهده است.



شکل ۲۸: نتیجه آموزش و ارزیابی با استفاده از SGD

همانطور که در تابع اتلاف مشاهده می شود، با هر بار آموزش مقدار اختلاف \hat{y} که همان خروجی مد نظر ما است از y_{train} که خروجی شبکه عصبی ما است، کمتر می شود و مقدار خطای ما کمتر می شود و به حدود صفر می رسد و نشان دهنده آن است که آموزش شبکه ما به درستی کار کرده است و رو به بهتر شدن می رود. هدف اصلی در این مسئله کمینه کردن مقدار تابع اتلاف است، به طوری که مدل توانایی خود را در پیش بینی یا تطبیق با داده های ورودی بهینه کند. در طول فرآیند آموزش، مدل بهبود می یابد و توانایی



پیش بینی بهتری را نسبت به داده های ورودی پیدا می کند. این به معنای این است که تابع اتلاف کاهش می یابد و مدل به سمت کمینه کردن خطا یا اختلاف بین خروجی مدل و مقادیر واقعی هدایت می شود. البته با توجه به ظاهر تابع اتلاف می شود این برداشت را کرد که شبکه امکان آموزش دیدن بیشتر را با افزایش تعداد ایپاک دارد اما به دلیل همگرا شدن و accuracy بالا نیازی به آموزش بیشتر وجود ندارد. تمام مراجل بالا را برای داده های ارزیابی نیز تکرار می کنیم با این تفاوت که این بار wها آپدیت نمی شوند و مقدار آخرین w بدست آمده در فرآیند آموزش را به عنوان ورودی برای داده های ارزیابی لحاظ می کنیم. برای محاسبه نتیجه ارزیابی از دو شاخصه accuracy که تابع محاسباتی آن را قبلا تعریف کردیم و همچنین f1-score از کتابخانه sklearn-metrics استفاده می کنیم.

```

1 x_test = np.hstack((np.ones((len(x_test), 1)), x_test))
2
3 y_hat_test = logistic_regression(x_test, w)
4
5 acc_test = accuracy(y_test, y_hat_test)
6
7 from sklearn.metrics import f1_score
8 f1score = f1_score(y_test, np.round(y_hat_test))
9
10 print("test accuracy:", acc_test)
11 print("test f1score:", f1score)

```

Code 33: import dataset and library

نمودار تابع اتلاف که در طول زمان آموزش شبکه عصبی رسم می شود، اطلاعات مفیدی را ارائه می دهد اما تنها از روی آن نمی توان به طور کامل و با قطعیت ظر دقیقی در مورد عملکرد نهایی مدل ارائه داد. دلایل آن هم می تواند نوسان داشتن نمودار تابع اتلاف در مراحل ابتدایی آموزش که به دلیل فرآیند آموزش و تنظیم پارامترهای مدل است و ممکن است در ادامه بهبود یابند، باشد. همچنین نمودار تابع اتلاف معمولا فقط نمایانگر عملکرد مدل بر روی داده های آموزشی است و اطلاعاتی درباره ی عملکرد واقعی مدل بر روی داده های جدید یا داده هایی که مدل آن ها را ندیده است، فراهم نمی کند.

اگر بخواهیم راه حلی ارائه دهیم برای اینکه بتوانیم پیش از مرحله ارزیابی در مورد عملکرد شبکه اظهار نظر کنیم، می توانیم داده ها آموزشی را به دو قسمت آموزش و اعتبارسنجی تقسیم کنیم. نحوه استفاده از این دسته جدید از داده هم به این صورت است که در هر ایپاک آموزشی که روی داده های آموزشی شبکه در حال یادگیری است، داده های اعتبارسنجی با وزن های آموزش دیده در همان ایپاک مورد ارزیابی قرار می گیرد و این فرآیند در کل طول همه ایپاک ها انجام می شود. مقایسه نمودار آموزش و اعتبارسنجی می تواند اطلاعات بسیار مفیدی از روند آموزش و نحوه یادگیری شبکه بدهد.



۴.۲ فرآیند آموزش و ارزیابی را با استفاده از یک طبقه بند خطی آماده پایتون (در `model-linear.sklearn`) (انجام داده و نتایج را مقایسه کنید. در حالت استفاده از دستورات آماده سایکیت لرن، آیا راهی برای نمایش نمودار تابع اتلاف وجود دارد؟ پیاده سازی کنید.

در این حالت از از طبقه بند آماده `LogisticRegression` استفاده می‌کنیم. با توجه به نتایج مناسبی که در مرحله عدم استفاده از کتابخانه‌های پایتونی به دست آمده بود، از پارامترهای پیش فرض استفاده می‌کنیم. کد این قسمت بصورت زیر است.

```
1 from sklearn.linear_model import LogisticRegression
2 y_train = y_train.ravel()
3 LR_model = LogisticRegression(penalty='l2',
4                               dual=False,
5                               tol = 1e-05,
6                               solver='sag',
7                               max_iter=1000,
8                               C=1,
9                               random_state=4)
10 LR_model.fit(x_train, y_train)
11 LR_model.predict(x_test)
12 print("LR_model train score = ", LR_model.score(x_train, y_train))
13 print("LR_model test score = ", LR_model.score(x_test, y_test))
```

Code 34: import dataset and library

خروجی مورد نظر به این صورت قابل مشاهده است.

```
LR_model train score = 1.0
LR_model test score = 1.0
```

شکل ۲۹: نتیجه آموزش و ارزیابی با استفاده از SGD

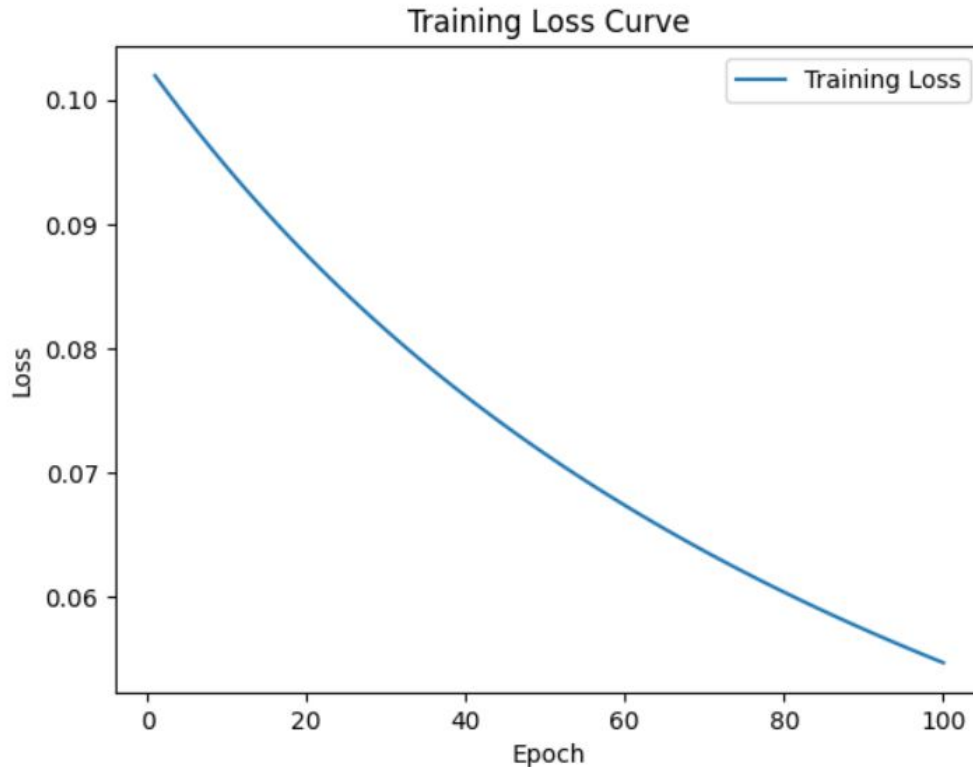
بله. امکان نمایش نمودار تابع اتلاف در حالت استفاده از دستورات آماده وجود دارد. در این حالت در ابتدا تمامی کتابخانه‌های مورد نیاز را در محیط `import` می‌کنیم و با استفاده از کلاس `SGDCeslassifier` مدل مورد نظرمان را روی داده‌ها پیاده می‌کنیم و در داخل `model` می‌ریزیم. در ادامه تابعی با نام `losses` تعریف می‌کنیم که در ادامه مقادیر اتلافی را در داخل آن بریزیم. در این قسمت برای بهبودی عملکرد مقدار ایپاک‌ها را زیاد و روی ۱۰۰۰ قرار می‌دهیم اما این مدل برای ایپاک‌های کمتر از این هم جوابگو خواهد بود. در ادامه ایپاک‌ها را به مدل مورد نظرمان اعمال می‌کنیم و تابع `losses` را `append` می‌کنیم.

```
1 from sklearn.metrics import log_loss
2 train_loss = []
```



```
3 # Train the classifier and collect loss values
4 epochs = 100
5 for epoch in range(epochs):
6     SGD_model.partial_fit(x_train, y_train, classes=np.unique(y_train))
7     loss = log_loss(y_train, SGD_model.predict_proba(x_train))
8     train_loss.append(loss)
9
10 # Plot the training loss curve
11 plt.plot(range(1, epochs + 1), train_loss, label='Training Loss')
12 plt.xlabel('Epoch')
13 plt.ylabel('Loss')
14 plt.title('Training Loss Curve')
15 plt.legend()
16 plt.show()
```

Code 35: import dataset and library

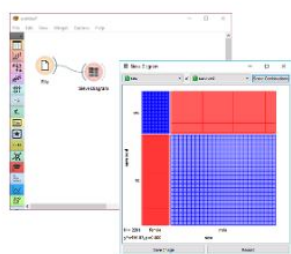


شکل ۳۰: نتیجه آموزش و ارزیابی با استفاده از SGD

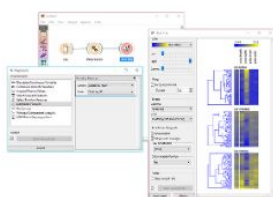


۵.۲ در مورد نرم افزار داده کاوی Orange و قابلیت های آن تحقیق کنید و سعی کنید این سوال یا یک مثال ساده تر را با استفاده از این نرم افزار پیاده سازی کنید (راهنمایی: می توانید از پیوندهای ۱ و ۲ و ۳ کمک بگیرید). پاسخ به این قسمت از سوال اختیاری و امتیازی است. می توانید عملکرد خود را به صورت تصویری و یا ویدیویی هم نشان دهید. مقدار نمره امتیازی، وابسته به جامعیت مثال بررسی شده و استفاده از ویژگی های مختلف این ابزار است.

این نرم افزار یک پلت فرم قدرتمند برای انجام تجزیه و تحلیل و تجسم داده ها، دیدن جریان داده ها و بهره وری بیشتر است. این یک پلت فرم منبع باز و تمیز و امکان افزودن قابلیت های بیشتر برای همه زمینه های علم را فراهم می کند. از جمله قابلیت های آن می توان به Heatmap visualisation, Receiver operating characteristics (ROC) analysis, Model-based feature scoring, CN2 rule induction اشاره کرد. در تصویر زیر به برخی دیگر از قابلیت های این نرم افزار اشاره شده است.



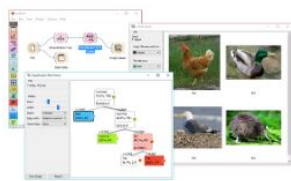
Sieve diagram on Titanic data set.



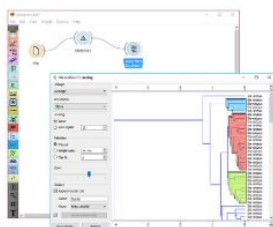
Heatmap visualisation.



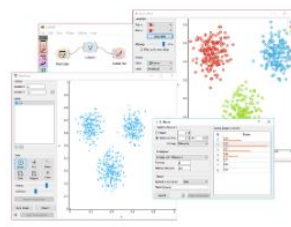
Explorative analysis with classification trees.



Data can contain references to images.



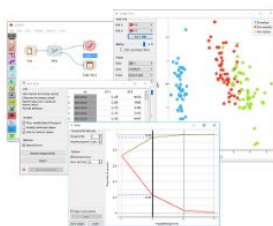
Hierarchical clustering supports interactive cluster selection.



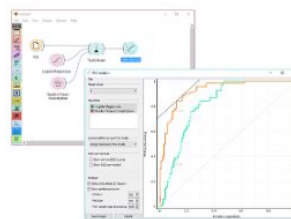
Playing with Paint Data and an automatic selection of clusters in k-Means.



Multidimensional scaling of Zoo data set reveals phylogeny groups.



Principal component analysis with scree diagram.



Receiver operating characteristics (ROC) analysis.

شکل ۳۱: نرم افزار Orange

ویدیوی عملکرد این نرم افزار در لینک زیر قابل دسترسی است.

https://drive.google.com/drive/folders/12pFGGxIkyYdCL5pybK_CpyKH4jNExdw?usp=sharing



۳ سوال سوم

یک دیتاست در زمینه آب و هوا با نام ۲۰۱۶-۲۰۰۶ Weather in Szeged را در نظر بگیرید. در این دیتاست هدف آن است که ارتباط بین Humidity با Temperature و هم چنین ارتباط بین Humidity و Apparent Temperature پیدا شده و با کمک داده های Humidity و Temperature تخمین انجام شود.

۱.۳ ابتدا هیت مپ ماتریس همبستگی و هیستوگرام پراکندگی ویژگی ها را رسم و تحلیل کنید.

ابتدا کتابخانه های مورد نیاز را فراخوانی می کنیم.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from matplotlib import pyplot as plt
6 import numpy as np
7 import matplotlib.pyplot as plt
8 from sklearn.model_selection import train_test_split
9 from sklearn.linear_model import LinearRegression
10 from sklearn.metrics import mean_squared_error
11 import statsmodels.api as sm
```

Code 36: import dataset and library

سپس فایل CSV دانلودی از مرجع مورد ذکر سوال را با دستور gdown فراخوانی می کنیم و آن را بصورت یک دیتافریم درمی آوریم.

```
1 %cd /content
2 !gdown 1fXFPECCGZ7-Kc8wXw8VhfGrsASlA0dAz
3
4 df = pd.read_csv("/content/weatherHistory.csv")
```

Code 37: import dataset and library

دیتافریم تشکیل شده بصورت زیر است.

هیت مپ ماتریس همبستگی که نشان دهنده میزان ارتباط داده های هر کدام از ویژگی ها است به صورت زیر قابل پیاده سازی است. در دیتاست مورد نظر با توجه به اینکه داده های برخی از ستون ها دارای مقادیر مورد نظر برای اعمال برای رسم این ماتریس بودند، ابتدا ستون های مورد نظر شناسایی شدند و برای اعمال انتخاب شدند.

```
1 # Exclude the last column from the DataFrame
```




	Formatted Date	Summary	Precip Type	Temperature (c)	Apparent Temperature (c)	Humidity	Wind Speed (km/h)	Wind Bearing (degrees)	Visibility (km)	Loud Cover	Pressure (millibars)	Daily Summary
0	2006-04-01 00:00:00.000 +0200	Partly Cloudy	rain	9.472222	7.388889	0.89	14.1197	251.0	15.8263	0.0	1015.13	Partly cloudy throughout the day.
1	2006-04-01 01:00:00.000 +0200	Partly Cloudy	rain	9.355556	7.227778	0.86	14.2646	259.0	15.8263	0.0	1015.63	Partly cloudy throughout the day.
2	2006-04-01 02:00:00.000 +0200	Mostly Cloudy	rain	9.377778	9.377778	0.89	3.9284	204.0	14.9569	0.0	1015.94	Partly cloudy throughout the day.
3	2006-04-01 03:00:00.000 +0200	Partly Cloudy	rain	8.288889	5.944444	0.83	14.1036	269.0	15.8263	0.0	1016.41	Partly cloudy throughout the day.
4	2006-04-01 04:00:00.000 +0200	Mostly Cloudy	rain	8.755556	6.977778	0.83	11.0446	259.0	15.8263	0.0	1016.51	Partly cloudy throughout the day.
...

شکل ۳۲: نتیجه آموزش و ارزیابی با استفاده از SGD

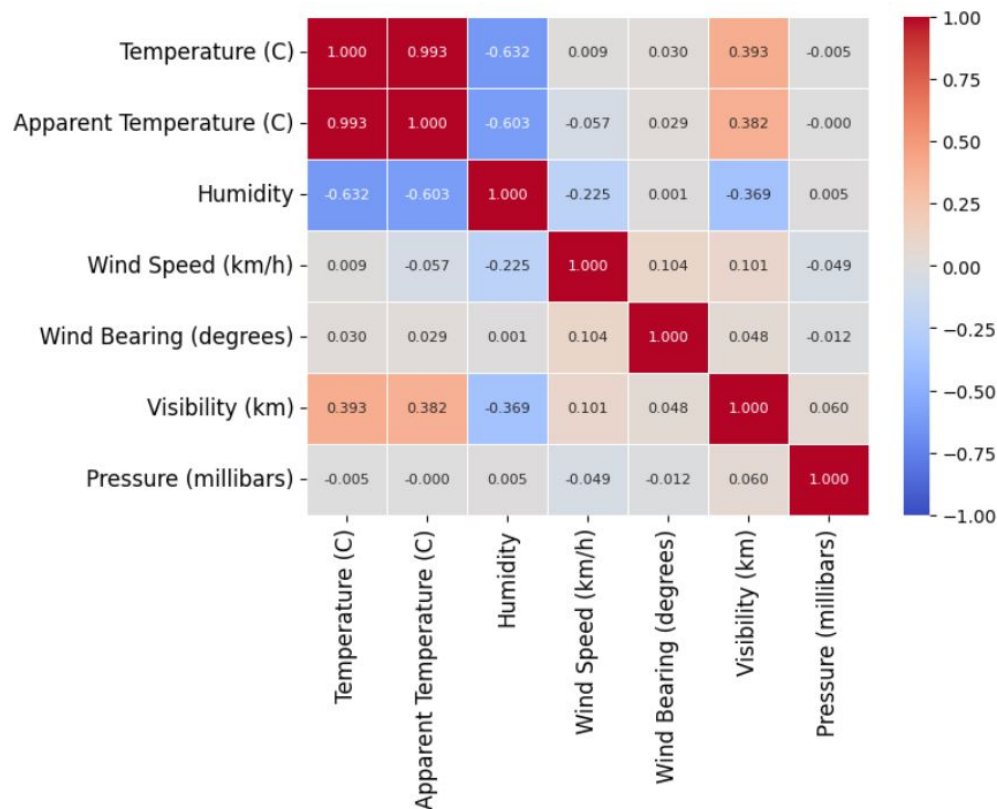
```

2 feature_df = df.iloc[:, [3,4,5,6,7,8,10]]
3
4 # Calculate correlation matrix
5 corr_matrix = feature_df.corr()
6
7 # Create heatmap using seaborn with a colormap from -1 to 1
8 plt.figure(figsize=(10, 10)) # Adjust the figure size as needed
9 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', linewidths=0.5,
10             annot_kws={"size": 8}, fmt='.3f',
11             yticklabels=corr_matrix.columns, vmin=-1, vmax=1) # Set vmin and
12             vmax
13
14 # Adjust font size of annotations
15 plt.xticks(fontsize=12)
16 plt.yticks(fontsize=12)
17
18 # Adjust margins of PDF file
19 plt.savefig('PIcS1.pdf', bbox_inches='tight')

```

Code 38: import dataset and library

ماتریس همبستگی به صورت زیر است. همانطور که قابل مشاهده است، ویژگی‌های Temperature، Apparent-Temperature، Humidity و Visibility نیز تا حدودی ارتباط بیشتری نسبت به سایر ویژگی‌ها با این دو ویژگی دارند. همچنین دو ویژگی Hu-



شکل ۳۳: نتیجه آموزش و ارزیابی با استفاده از SGD

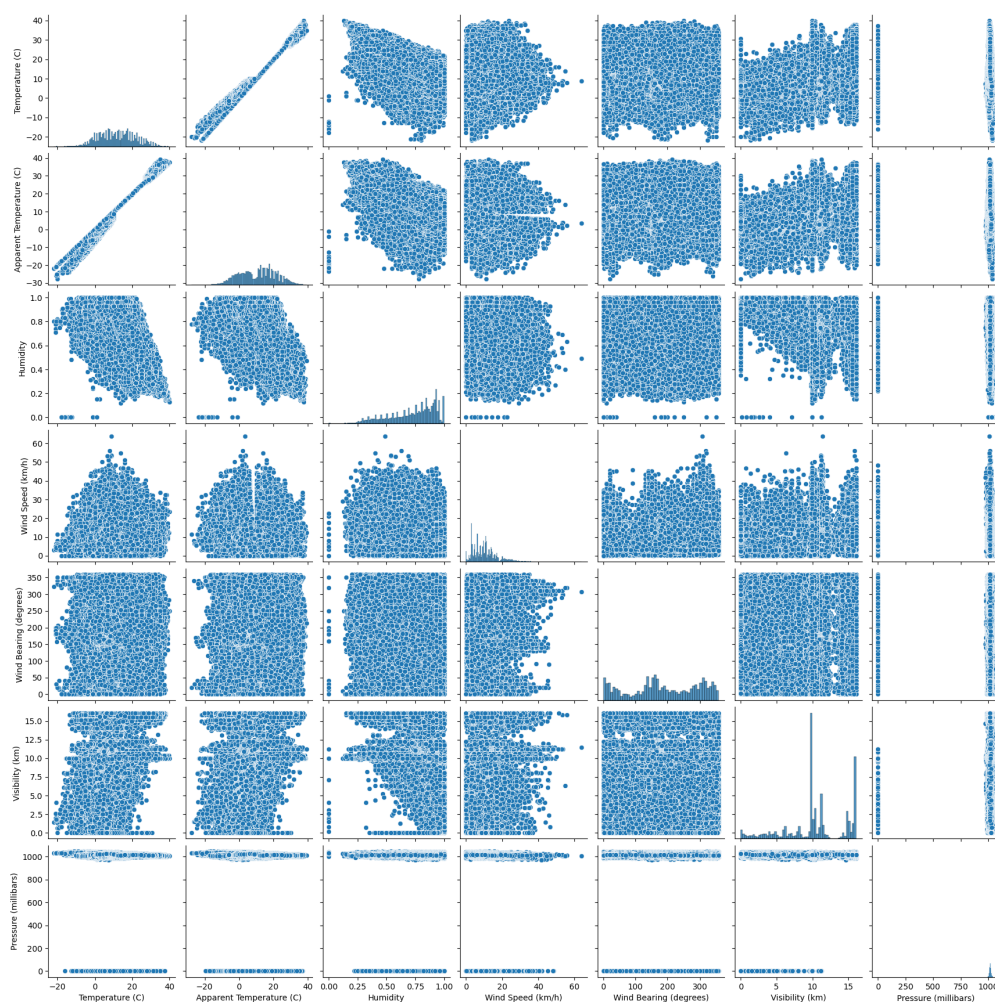
همچنین نمودار پراکندگی ویژگی‌های مختلف نسبت به یکدیگر بصورت زیر قابل مشاهده است. در صورت سوال تاکید شده است که از داده‌های مربوط به ویژگی‌های Temperature Humidity و Apparent-Temperature استفاده شود اما با توجه به این که ویژگی Visibility نیز تا حدودی همبستگی مناسبی دارد، فعلا در ادامه کار روی این ویژگی نیز بررسی‌هایی را صورت می‌دهیم تا نتیجه مناسبی رقم بخورد. با استفاده از این کد این ویژگی‌ها را بصورت جداگانه استخراج می‌کنیم.

```
1 Humidity = df['Humidity'].values
2 Temperature = df["Temperature (C)"].values
3 Apparent_Temperature = df["Apparent Temperature (C)"].values
4 Visibility = df["Visibility (km)"].values
```

Code 39: import dataset and library

در قسمت بعد ویژگی‌هایی که همبستگی بیشتری داشتند را از نظر هیستوگرام پراکندگی ویژگی‌ها بررسی می‌کنیم. کد زیر این هیستوگرام‌ها را رسم می‌کند. قابل مشاهده است که ویژگی‌های Temperature و Apparent-Temperature دارای پراکندگی ویژگی بسیار مشابه یکدیگراند. ویژگی Humidity دارای پراکندگی بیشتر در مقادیر نزدیک به یک است و ویژگی Visibility دارای پراکندگی بیشتری در حد متوسط آن است.

```
1 # Plot histogram for 'Temperature (C)'
```



شکل ۳۴: نتیجه آموزش و ارزیابی با استفاده از SGD

```
2 plt.figure(figsize=(12, 2.3))
3 plt.subplot(1, 4, 1)
4 feature_df['Temperature (C)'].plot(kind='hist', bins=20, title='Temperature')
5 plt.gca().spines[['top', 'right']].set_visible(False)
6
7 # Plot histogram for 'Apparent Temperature (C)'
8 plt.subplot(1, 4, 2)
9 feature_df['Apparent Temperature (C)'].plot(kind='hist', bins=20, title='
    Apparent Temperature')
10 plt.gca().spines[['top', 'right']].set_visible(False)
11
12 # Plot histogram for 'Humidity'
```

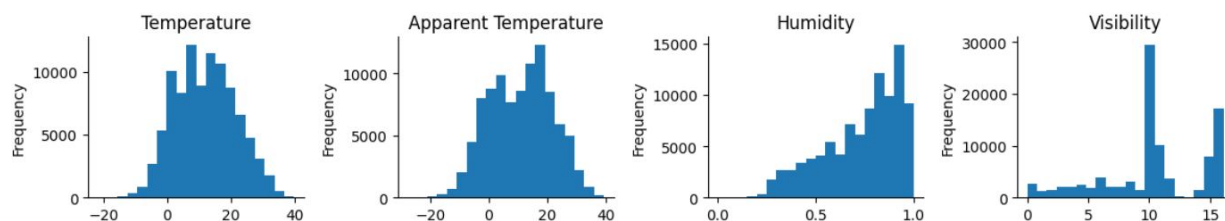


```

13 plt.subplot(1, 4, 3)
14 feature_df['Humidity'].plot(kind='hist', bins=20, title='Humidity')
15 plt.gca().spines[['top', 'right']].set_visible(False)
16
17 # Plot histogram for 'Visibility'
18 plt.subplot(1, 4, 4)
19 feature_df['Visibility (km)'].plot(kind='hist', bins=20, title='Visibility')
20 plt.gca().spines[['top', 'right']].set_visible(False)
21
22 # Adjust layout
23 plt.tight_layout()
24
25 # Show plot
26 plt.show()

```

Code 40: import dataset and library



شکل ۳۵: نتیجه آموزش و ارزیابی با استفاده از SGD

۲.۳ روی این دیتاست، تخمین LS و RLS را با تنظیم پارامترهای مناسب اعمال کنید. نتایج به دست آمده را با محاسبه خطاها و رسم نمودارهای مناسب برای هر دو مدل با هم مقایسه و تحلیل کنید.

در صورت مسئله عنوان شده است که با کمک داده‌های Humidity، Temperature تخمین صورت بپذیرد. بنابراین در این قسمت از این دو ویژگی به عنوان ویژگی‌هایی که شبکه RL به عنوان ورودی می‌پذیرد استفاده می‌کنیم و ویژگی Apparent-Temperature را به عنوان خروجی که شبکه سعی در تخمین زدن آن دارد انتخاب می‌کنیم. ابتدا شبکه RL مورد نظر را به صورت زیر فراخوانی می‌کنیم. در این قسمت از کد مورد استفاده در کلاس حل تمرین در این قسمت استفاده شده است.

```

1 class LinearRegressionLS:
2     def __init__(self):

```



```
3     self.coefficients = None
4
5     def fit(self, X, y):
6         # Add a column of ones to account for the intercept term
7         X = np.column_stack((np.ones(len(X)), X))
8
9         # Compute the coefficients using the least squares method
10        self.coefficients = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(y)
11
12    def predict(self, X):
13        # Add a column of ones to account for the intercept term
14        X = np.column_stack((np.ones(len(X)), X))
15
16        # Predict the target variable
17        return X.dot(self.coefficients)
```

Code 41: import dataset and library

از ویژگی‌های Humidity و Temperature و Visibility مطابق صورت مسئله استفاده شده است و سعی می‌شود مدل بتواند تخمین مناسبی از Apparent-Temperature پیش‌بینی شود. سپس داده‌های فراهم شده به دو قسمت آموزش و ارزیابی تقسیم می‌شوند و از مدلی که پیش‌تر توضیح داده شد استفاده شده و پیش‌بینی و مقدار خطای MSE محاسبه می‌شود.

```
1 X = np.concatenate((Humidity, Temperature, Visibility), axis=1)
2 y = Apparent_Temperature
3
4 # Split the dataset into training and testing sets
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
6             random_state=4)
7 print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
8 print()
9
10 # Initialize and fit the linear regression model using least squares
11 model = LinearRegressionLS()
12 model.fit(X_train, y_train)
13
14 # Make predictions on the testing set
```



```

15 y_pred = model.predict(X_test)
16
17 # Calculate Mean Squared Error
18 mse = mean_squared_error(y_test, y_pred)
19 print("Mean Squared Error:", mse)

```

Code 42: import dataset and library

خروجی نهایی به صورت زیر است. مقدار خطای MSE نشان‌دهنده میزان نزدیکی مقادیر پیش‌بینی شده و مقدار ویژگی هدف است که هر چه این خطا کمتر باشد، یعنی مدل عملکرد بهتری در تخمین داشته است.

```
(77162, 3) (19291, 3) (77162, 1) (19291, 1)
```

```
Mean Squared Error: 1.5746506582916415
```

شکل ۳۶: نتیجه آموزش و ارزیابی با استفاده از SGD

به صورت زیر می‌توان نمودار مقادیر تخمینی شبکه را بر حسب خروجی هدف که قصد تخمین زدن آن را داشتیم رسم کرد. این نمودار حاوی مقادیر تمامی داده‌های مقادیر واقعی و تخمینی و خطی به عنوان خط تخمینی است.

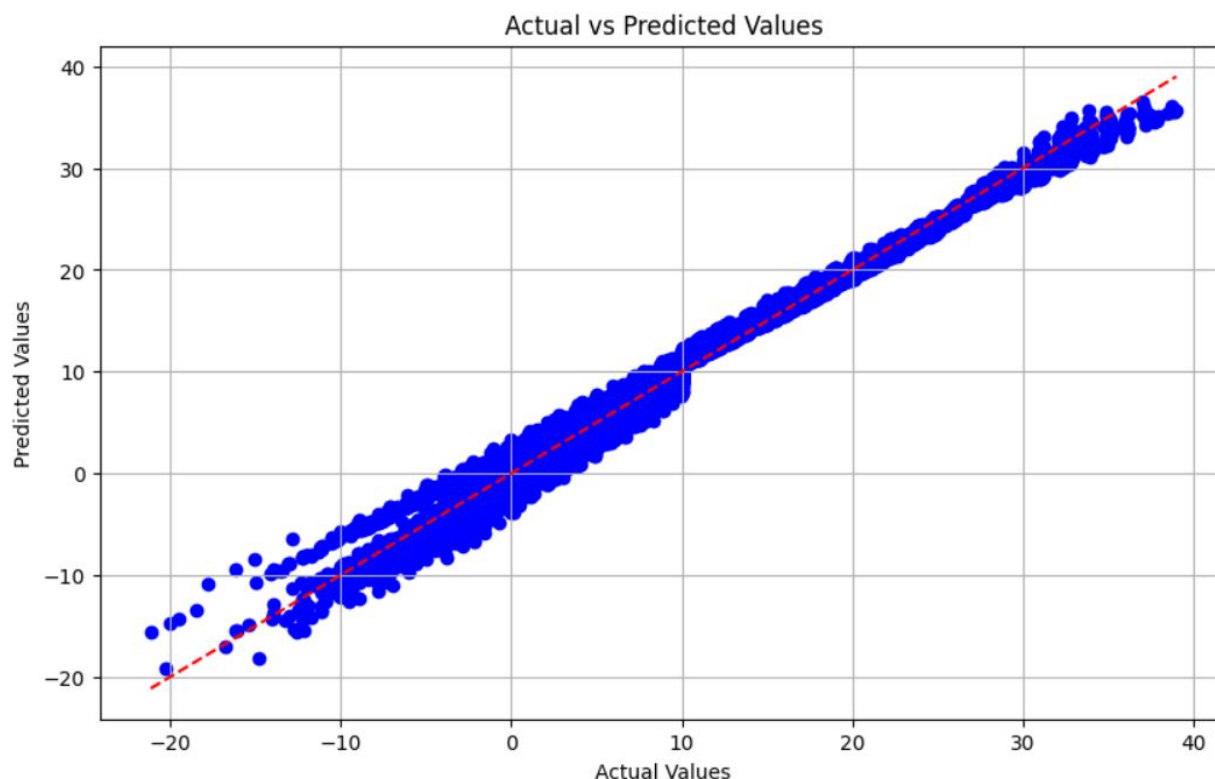
```

1 # Plot actual vs predicted values
2 plt.figure(figsize=(10, 6))
3 plt.scatter(y_test, y_pred, color='blue')
4 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='
    red', linestyle='--')
5 plt.xlabel('Actual Values')
6 plt.ylabel('Predicted Values')
7 plt.title('Actual vs Predicted Values')
8 plt.grid(True)
9 plt.show()

```

Code 43: import dataset and library

برای روش RLS از کد زیر بهره می‌بریم. در این حالت یک کلاس تشکیل شده است که یک ورودی با عنوان forgetting-factor می‌گیرد. این ورودی با میزان به یادآوری مقادیر قبلی ارتباط دارد که هرچی این مقدار بیشتر باشد، به مقادیر جدیدتری که اخیراً شبکه تولید کرده است وابسته‌تر می‌شود. بهترین مقدار این پارامتر مطابق آزمایش‌ها برابر مقدار ۹۹۰۰ یا ۹۰۰ است که ما هم در این قسمت از همین مقدارهای اسفاده کرده‌ایم اما در ادامه صحت این ادعا را بررسی می‌کنیم.



شکل ۳۷: نتیجه آموزش و ارزیابی با استفاده از SGD

```
1 class RecursiveLeastSquares:
2     def __init__(self, n_features, forgetting_factor=0.99):
3         self.n_features = n_features
4         self.forgetting_factor = forgetting_factor
5         self.theta = np.zeros((n_features, 1)) # Initialize model parameters
6         self.P = np.eye(n_features) # Initialize covariance matrix
7
8     def fit(self, X, y):
9         errors = []
10        for i in range(len(X)):
11            x_i = X[i].reshape(-1, 1)
12            y_i = y[i]
13
14            # Predict
15            y_pred = np.dot(x_i.T, self.theta)
16
17            # Update
```



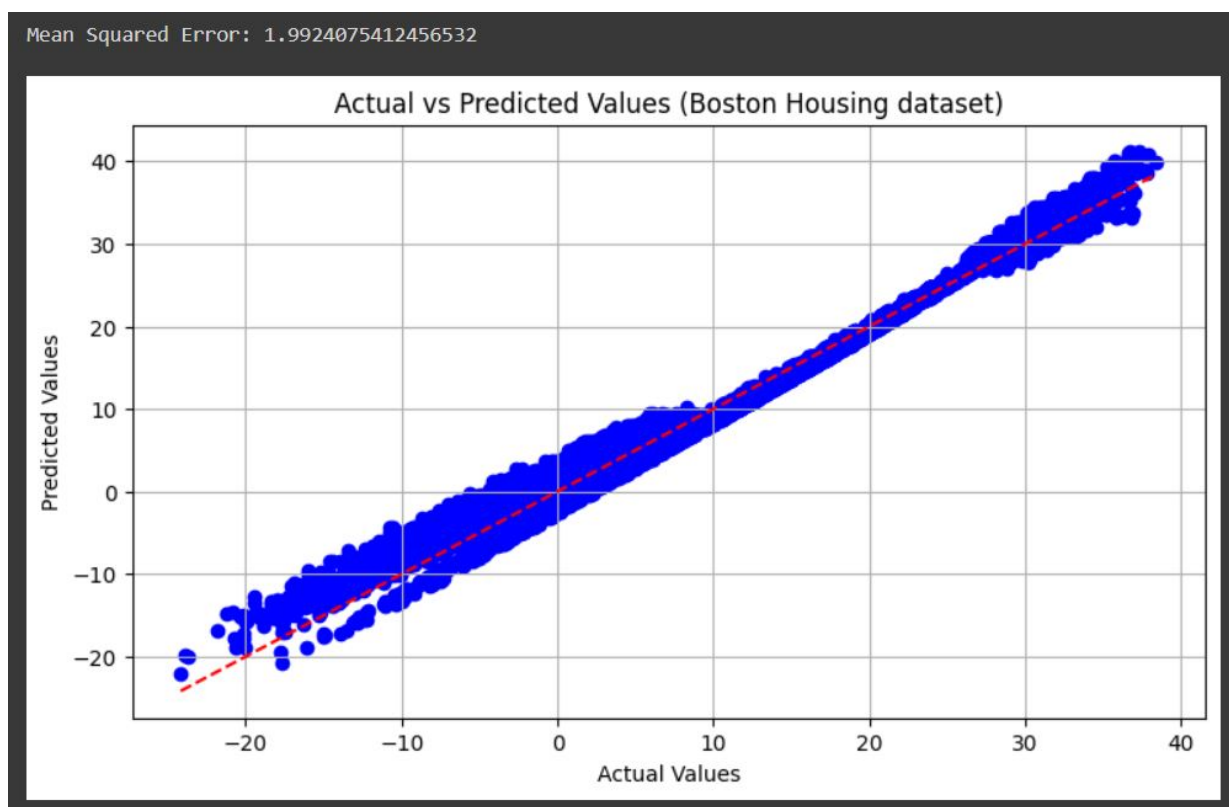
```
18         error = y_i - y_pred
19         errors.append(error)
20         K = np.dot(self.P, x_i) / (self.forgetting_factor + np.dot(np.dot(
x_i.T, self.P), x_i))
21         self.theta = self.theta + np.dot(K, error)
22         self.P = (1 / self.forgetting_factor) * (self.P - np.dot(K, np.dot
(x_i.T, self.P)))
23
24     return errors
25
26     def predict(self, X):
27         return np.dot(X, self.theta)
28
29 X = np.concatenate((Humidity, Visibility, Temperature), axis=1)
30 y = Apparent_Temperature
31
32 # Split the dataset into training and testing sets
33 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
34
35 # Initialize and fit the RLS model
36 rls = RecursiveLeastSquares(n_features=X_train.shape[1], forgetting_factor
=0.99)
37 errors = rls.fit(X_train, y_train)
38
39 # Make predictions
40 y_pred = rls.predict(X_test)
41
42 # Calculate Mean Squared Error
43 mse = np.mean(np.array(errors)**2)
44 print("Mean Squared Error:", mse)
45
46 # Plot actual vs predicted values
47 plt.figure(figsize=(10, 6))
48 plt.scatter(y_test, y_pred, color='blue')
49 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color=''
```




```
    red', linestyle='--')
50 plt.xlabel('Actual Values')
51 plt.ylabel('Predicted Values')
52 plt.title('Actual vs Predicted Values (Boston Housing dataset)')
53 plt.grid(True)
54 plt.show()
```

Code 44: import dataset and library

مقدار خطای MSE و نمودار مقدار تخمینی توسط شبکه بر حسب مقادیر واقعی در تصویر زیر قابل مشاهده است. مقدار خطا به مقادیر کم رسیده است و شبکه توانسته است خطی که به طور مناسبی بر مقادیر تخمینی و واقعی قرار می‌گیرد را شناسایی و رسم کند.



شکل ۳۸: نتیجه آموزش و ارزیابی با استفاده از SGD

اما برای بررسی صحت ادعای بهترین مقدار forgetting-factor در یک حلقهٔ یکرار مقادیر مختلف این پارامتر را مورد آزمایش قرار می‌دهیم.

```
1 mse_list = []
2 for ff in np.linspace(0.5, 0.99, 30):
3     # Initialize and fit the RLS model
```



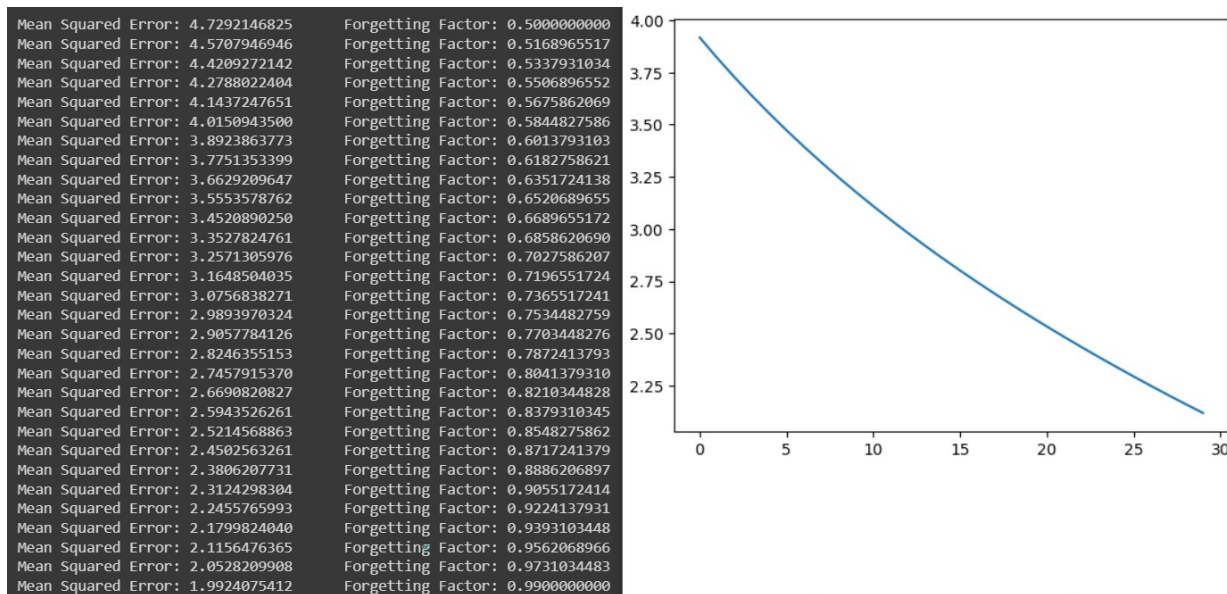
```

4 rls = RecursiveLeastSquares(n_features=X_train.shape[1], forgetting_factor
    = ff)
5 errors = rls.fit(X_train, y_train)
6
7 # Make predictions
8 y_pred = rls.predict(X_test)
9
10 # Calculate Mean Squared Error
11 mse = np.mean(np.array(errors)**2)
12 mse_list.append(mse)
13 # print("Mean Squared Error:", mse, "           Forgetting Factor:", ff)
14 print("Mean Squared Error: {:.10f}           Forgetting Factor: {:.10f}".format
    (mse, ff))

```

Code 45: import dataset and library

نتیجه به صورت زیر قابل مشاهده است که هر چه به سمت مقادیر نزدیک ۱ می‌رویم نتیجه بهبود می‌یابد.



شکل ۳۹: نتیجه آموزش و ارزیابی با استفاده از SGD



۳.۳ در مورد Weighted Least Square توضیح دهید و آن را روی دیتاست داده شده اعمال کنید.

حداقل مربعات وزن دار یک تکنیک رگرسیونی است که با تخصیص وزن به مشاهدات بر اساس واریانس های تخمین زده شده، به بررسی ناهمسانی می پردازد و امکان برآورد پارامتر قوی تر و عملکرد مدل را بهبود می بخشد. به عبارتی دیگر حداقل مربعات وزنی (WLS) نوعی تغییر از روش حداقل مربعات معمولی (OLS) است که در تحلیل رگرسیون خطی استفاده می شود. در OLS، هر نقطه داده از نظر تأثیر آن بر تخمین ضرایب رگرسیون به طور مساوی رفتار می شود. با این حال، در برخی موارد، تغییرپذیری نقاط داده ممکن است برابر نباشد و برخی از مشاهدات ممکن است قابل اعتمادتر باشند یا واریانس کمتری نسبت به سایرین داشته باشند.

حداقل مربعات وزنی با تخصیص وزن به هر نقطه داده بر اساس واریانس یا پایایی آنها، به این موضوع می پردازد. ایده اصلی این است که به مشاهداتی که قابل اعتمادتر در نظر گرفته می شوند یا دارای تنوع کمتری هستند و وزن کمتری به مشاهداتی که کمتر قابل اعتماد تلقی می شوند یا دارای تنوع بالاتری هستند، وزن بیشتری بدهیم. این کار با اعمال ماتریس W به فرمول محاسباتی LS صورت می پذیرد که به شکل تصویر زیر اعمال می شود.

$$W = \begin{pmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_n \end{pmatrix}$$

$$\begin{aligned} \hat{\beta}_{WLS} &= \arg \min_{\beta} \sum_{i=1}^n \epsilon_i^{*2} \\ &= (X^T W X)^{-1} X^T W Y \end{aligned}$$

شکل ۴۰: نتیجه آموزش و ارزیابی با استفاده از SGD

برای اعمال این روش بر روی مجموعه داده مورد استفاده در این سوال از کد زیر استفاده می کنیم. در این حالت با تعریف ماتریس وزن از حالت LS به حالت WLS وارد می شویم و به این طریق وزن های مورد نظر را اعمال می کنیم.

```
1 X = np.concatenate((Humidity, Temperature, Visibility), axis=1)
2 y = Apparent_Temperature
3
4 # Split the dataset into training and testing sets
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
6                                                    random_state=4)
7 print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
8 print()
```



```

9 error_variance = 2 # Modify this value based on your estimation
10
11 # Calculate weights based on the estimated variance
12 weights = 1 / error_variance
13
14 # Fit the weighted least squares model
15 X_with_intercept = sm.add_constant(X_train) # Add intercept term
16 model = sm.WLS(y_train, X_with_intercept, weights=weights)
17 result = model.fit()
18
19 # Print the model summary
20 print(result.summary())

```

Code 46: import dataset and library

خلاصه‌ای از پارامترهای شبکه آموزش داده شده به صورت زیر است و همانطور که قابل مشاهده است مقادیر مانند R-squared و Adj. R-squared که به مقادیر نزدیک ۱ رسیده‌اند نشان‌دهنده همگرا شدن الگوریتم مورد نظر به تخمین مناسب است.

(77162, 3) (19291, 3) (77162, 1) (19291, 1)					
WLS Regression Results					
Dep. Variable:	y	R-squared:	0.986		
Model:	WLS	Adj. R-squared:	0.986		
Method:	Least Squares	F-statistic:	1.866e+06		
Date:	Fri, 12 Apr 2024	Prob (F-statistic):	0.00		
Time:	23:26:51	Log-likelihood:	-1.2655e+05		
No. Observations:	77162	AIC:	2.531e+05		
Df Residuals:	77158	BIC:	2.531e+05		
df Model:	3				
covariance Type:	nonrobust				
	coef	std err	t	P> t	[0.025 0.975]
const	-4.3176	0.031	-139.838	0.000	-4.378 -4.257
x1	2.2417	0.039	74.512	0.000	2.183 2.301
x2	1.1423	0.001	1833.235	0.000	1.141 1.143
x3	-0.0103	0.001	-8.686	0.000	-0.013 -0.008
Omnibus:	1603.053	Durbin-Watson:			1.981
prob(Omnibus):	0.000	Jarque-Bera (JB):			2739.820
skew:	-0.183	Prob(JB):			0.00
kurtosis:	3.847	Cond. No.			171.
Notes:					
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.					

شکل ۴۱: نتیجه آموزش و ارزیابی با استفاده از SGD

برای اطمینان از صحت روش مورد استفاده، نیاز است که مشابه موارد قبلی نمودار تخمین بر اساس مقادیر واقعی رسم شود. قابل مشاهده است که تخمین مانند موارد قبلی همگرا شده و توانسته به طور مناسبی خطی منطبق و تخمینی بر داده‌های تخمینی بر اساس مقادیر واقعی رسم کند.

```

1 X_test_with_intercept = sm.add_constant(X_test)
2 y_pred = result.predict(X_test_with_intercept)
3
4 # Plot Predicted vs. Actual Values
5 plt.figure(figsize=(8, 6))
6 plt.scatter(y_pred, y_test, color='blue', alpha=0.5)

```

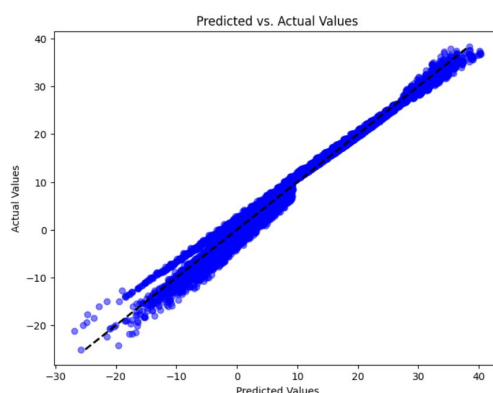


```

7 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw
          =2) # Diagonal line for reference
8 plt.xlabel('Predicted Values')
9 plt.ylabel('Actual Values')
10 plt.title('Predicted vs. Actual Values')
11 plt.show()

```

Code 47: import dataset and library



شکل ۴۲: نتیجه آموزش و ارزیابی با استفاده از SGD

۴.۳ در مورد الگوریتم **QR-Decomposition-Based RLS** تحقیق کنید. پاسخ به این قسمت از سوال اختیاری و امتیازی است.

استفاده از تجزیه QR برای مثلی کردن ماتریس داده های ورودی منجر به یک روش جایگزین برای اجرای روش حداقل مربعات بازگشتی (RLS) می شود. مزایای اصلی ناشی از الگوریتم حداقل مربعات بازگشتی مبتنی بر تجزیه QR، اجرای احتمالی آن در آرایه های سیستمیک و بهبود رفتار عددی آن با در نظر گرفتن اثرات کوانتیزاسیون است.

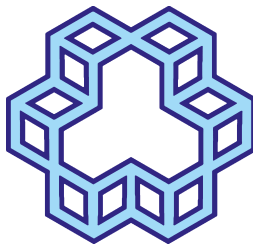
به عبارت دیگر حداقل مربعات بازگشتی مبتنی بر تجزیه روشی است که برای تخمین پارامتر در مسائل فیلتر تطبیقی و حداقل مربعات بازگشتی استفاده می شود. این تکنیک تجزیه QR را با به روزرسانی بازگشتی ترکیب می کند تا با در دسترس قرار گرفتن داده های جدید، راه حل حداقل مربعات را به طور مؤثر محاسبه کند.

الگوریتم RLS بر پایه Decomposition QR از ترکیب دو تکنیک استفاده می کند: Least Recursive و Decomposition QR. این الگوریتم برای پردازش سیگنال های پویا و تخمین پارامترها استفاده می شود. در این الگوریتم، ابتدا ماتریس مورد نظر را با استفاده از روش Decomposition QR به دو ماتریس قاری (یا orthogonal) و ماتریس سه گانه بالایی (یا triangular) upper Squares. سپس از روش Squares Least Recursive برای به روزرسانی پاسخ به صورت بازگشتی با هر آمدن داده جدید استفاده می کنیم، بدون نیاز به محاسبه مجدد کل پاسخ از ابتدا. این ترکیب از دو روش، باعث افزایش کارایی و کاهش پیچیدگی محاسباتی مسئله می شود، به ویژه زمانی که با مجموعه داده های بزرگ یا نیاز به پردازش به صورت زمان بندی شده روبرو هستیم.



۴ عنوان سوال چهارم

در این قسمت با نحوه درج اشکال آشنا می شوید:



شکل ۴۳: شکل شماره ۱

۵ عنوان سوال پنجم

در این قسمت با نحوه درج جداول آشنا می شوید:

جدول ۱: جدول شماره ۱

خانه شماره ۱	خانه شماره ۲	خانه شماره ۳
خانه شماره ۴	خانه شماره ۵	خانه شماره ۶
خانه شماره ۷	خانه شماره ۸	خانه شماره ۹

۶ عنوان سوال ششم

در این قسمت با نحوه درج انواع لیست ها آشنا می شوید:

۱.۶ عنوان بخش اول سوال ششم

- مورد اول
- مورد دوم

۲.۶ عنوان بخش دوم سوال ششم

۱. مورد شماره ۱
۲. مورد شماره ۲



۷ عنوان سوال هفتم

در این قسمت با نحوه درج برنامه‌ها آشنا می‌شوید:

```
1 # This program prints Hello, world!
2
3 print('Hello, world!')
```

Code 48: My Caption (Python)

```
1 clc; clear all; close all;
2 disp('Hello world!')
```

Code 49: My Caption (MATLAB)

```
1 // Your First C++ Program
2
3 #include <iostream>
4
5 int main() {
6     std::cout << "Hello World!";
7     return 0;
8 }
```

Code 50: My Caption (C++)

```
1 #include <stdio.h>
2 int main()
3 {
4     printf("Hello world!\n");
5     return 0;
6 }
```

Code 51: My Caption (C)



۸ عنوان سوال هشتم

در این قسمت با نحوه ارجاع دادن آشنا می شوید.

۱.۸ عنوان بخش اول سوال هشتم

در این قسمت با نحوه ارجاع به سایر منابع آشنا می شوید:

به صفحه درس تشخیص و شناسایی عیب ارجاع داده می شود [b۱]. به این کتاب ها ارجاع داده می شود [b۲][b۳]. برای وارد کردن ارجاع می توانید از انتهای فایل main.tex استفاده کنید و یا با تغییر قالب مرجع نویسی، به فایل bibliography.bib مراجعه کرده و فرمت bib را وارد کنید.

۲.۸ عنوان بخش دوم سوال هشتم

اگر می خواهید به یک شکل، جدول، یا بخش ارجاع دهید می توانید به دو صورتی که در ادامه آمده عمل کنید (حالت اول توصیه می شود):

۱. مورد شماره ۱: پاسخ سوال ۱،؟؟، شکل ۴۳، جدول ۱، ۴.

۲. مورد شماره ۲: سوال اول.

۳.۸ عنوان بخش سوم سوال هشتم

اگر می خواهید به یک پایگاه اینترنتی ارجاع دهید، می توانید از این دستور هم استفاده کنید: گیتهاب (GitHub).

۹ ضمیمه

برای آشنایی بیشتر با \LaTeX ، با جست و جو در اینترنت منابع مفیدی خواهید یافت.

منابع

[۱] صفحه 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. Koutroubas, "Pattern recognition", Fourth Edition, Academic Press, 2009.