

مخزن گیت‌هاب

برای مشاهده و دانلود کدهای مربوط به این پروژه، می‌توانید به مخزن گیت‌هاب مراجعه کنید:

مخزن گیت‌هاب پروژه‌ها

دفترچه گوگل کولب

همچنین دفترچه‌ی گوگل کولب مربوط به تمرین را از لینک زیر مشاهده کنید:

Colab Google در AI۴۰۳۲_HW۱

گزارش تمرین ۱

هوش مصنوعی

دکتر علیاری

سینا حسن پور شماره دانشجویی: ۴۰۲۱۶۷۲۳

۱۸ اسفند ۱۴۰۳

۱ مقدار جبر خطی

I ضرب و ابعاد

ماتریس‌های زیر را در نظر بگیرید:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \\ b_{41} & b_{42} \end{bmatrix}$$

(آ) طبق قاعده‌ی ضرب ماتریس‌ها:

اگر ماتریس A ابعاد $m \times n$ داشته باشد و ماتریس B ابعاد $n \times p$ داشته باشد، می‌توانیم AB را حساب کنیم و ماتریس حاصل ابعاد $m \times p$ دارد. اگر تعداد ستون‌های A با تعداد سطرهای B برابر نباشد، نمی‌توانیم ضرب را انجام بدهیم! ابعاد این ماتریس‌ها:

• ابعاد A : 2×3

• ابعاد B : 4×2

• ابعاد B^T (ترانهایه‌ی B) : 2×4

• ابعاد A^T (ترانهایه‌ی A) : 3×2

بررسی امکان‌پذیری ضرب‌ها با بررسی ابعاد، می‌توان دریافت که: - ضرب BA امکان‌پذیر است و نتیجه‌ی آن یک ماتریس با ابعاد 4×3 خواهد بود. - ترانهایه‌ی B ، یعنی B^T ، دارای ابعاد 2×4 است. - ضرب‌های $A^T B$ و $B^T A$ به دلیل اینکه برخلاف قاعده ضرب بودند، تعریف نشده‌اند.

ب) محاسبه‌ی BA و B^T

محاسبه‌ی B^T

ترانهاده‌ی ماتریس B به صورت زیر است:

$$B^T = \begin{bmatrix} b_{11} & b_{21} & b_{31} & b_{41} \\ b_{12} & b_{22} & b_{32} & b_{42} \end{bmatrix}$$

محاسبه‌ی BA

عملیات ضرب ماتریسی برای BA انجام می‌شود:

$$BA = \begin{bmatrix} (b_{11}a_{11} + b_{12}a_{21}) & (b_{11}a_{12} + b_{12}a_{22}) & (b_{11}a_{13} + b_{12}a_{23}) \\ (b_{21}a_{11} + b_{22}a_{21}) & (b_{21}a_{12} + b_{22}a_{22}) & (b_{21}a_{13} + b_{22}a_{23}) \\ (b_{31}a_{11} + b_{32}a_{21}) & (b_{31}a_{12} + b_{32}a_{22}) & (b_{31}a_{13} + b_{32}a_{23}) \\ (b_{41}a_{11} + b_{42}a_{21}) & (b_{41}a_{12} + b_{42}a_{22}) & (b_{41}a_{13} + b_{42}a_{23}) \end{bmatrix}$$

II ضرب و ابعاد، دشوارتر

پاسخ به سؤالات (الف) و (ب)

الف) بررسی ابعاد

• فرض کنید $x^{(i)} \in \mathbb{R}^{1 \times 2}$ ، یعنی $x^{(i)} = [x_1^{(i)}, x_2^{(i)}]$.

• اگر $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in \mathbb{R}^{2 \times 1}$ ، آنگاه ضرب $x^{(i)} \theta$ که ابعاد $(1 \times 2) \times (2 \times 1)$ دارد، یک اسکالر (1×1) خواهد بود:

$$x^{(i)} \theta = [x_1^{(i)}, x_2^{(i)}] \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = x_1^{(i)} \theta_1 + x_2^{(i)} \theta_2 \in \mathbb{R}.$$

• بعد θ در این مثال (2×1) است و بنابراین بعد $\theta^T x^{(i)}$ یا $x^{(i)} \theta$ اسکالر (1×1) است.

• وقتی همه‌ی نمونه‌ها $x^{(i)} \in \mathbb{R}^{1 \times 2}$ (برای $i = 1, \dots, n$) را در یک ماتریس جمع کنیم، ماتریس

$$X \in \mathbb{R}^{n \times 2}$$

به دست می‌آید که هر سطر آن یکی از بردارهای $x^{(i)}$ است:

$$X = \begin{bmatrix} -(x^{(1)})^- \\ -(x^{(2)})^- \\ \vdots \\ -(x^{(n)})^- \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots \\ x_1^{(n)} & x_2^{(n)} \end{bmatrix}.$$

بعد این ماتریس $n \times 2$ است.

• ضرب $X \theta$ آنگاه ابعاد $n \times 1$ خواهد داشت:

$$X \theta \in \mathbb{R}^{n \times 1}.$$

ب) بازنویسی تابع هزینه و محاسبه‌ی گرادیان تابع هزینه (Cost Function) رگرسیون خطی معمولاً به صورت

$$J(\theta) = \sum_{i=1}^n \left(x^{(i)} \theta - y^{(i)} \right)^2$$

تعریف می‌شود.

• نوشتن در قالب ضرب برداری:

اگر $X \in \mathbb{R}^{n \times 2}$, $\theta \in \mathbb{R}^{2 \times 1}$, $y \in \mathbb{R}^{n \times 1}$, آنگاه

$$J(\theta) = (X\theta - y)^T (X\theta - y).$$

این همان جمع مربعات مؤلفه‌های بردار $(X\theta - y)$ است، یعنی:

$$(X\theta - y)^T (X\theta - y) = \sum_{i=1}^n \left([x^{(i)} \theta] - y^{(i)} \right)^2.$$

• بسط و محاسبه‌ی گرادیان:

$$J(\theta) = (X\theta - y)^T (X\theta - y) = \theta^T X^T X \theta - 2y^T X \theta + y^T y.$$

مشتق این عبارت نسبت به θ (یا $\nabla_{\theta} J(\theta)$) برابر است با:

$$\nabla_{\theta} J(\theta) = 2 X^T (X\theta - y).$$

(اگر در تعریف اولیه‌ی $J(\theta)$ ضریب $\frac{1}{2}$ به کار رفته باشد، ضریب ۲ در گرادیان حذف می‌گردد.)

نتیجه‌گیری نهایی

- بعد θ در این مثال (2×1) است.
- بعد $x^{(i)}$ (1×2) و بنابراین $x^{(i)} \theta$ اسکالر (1×1) .
- با کنار هم گذاشتن تمامی $x^{(i)}$ ها، ماتریس X با بعد $(n \times 2)$ ساخته می‌شود و ضرب $X\theta$ در ابعاد $n \times 1$ خواهد بود.
- تابع هزینه $J(\theta)$ را می‌توان به صورت $(X\theta - y)^T (X\theta - y)$ نوشت و گرادیان آن نسبت به θ نیز $2 X^T (X\theta - y)$ به دست می‌آید.

۲ پردازش داده

Dataset CWRU I

(آ دریافت داده :

فرمت فایل داده شده 'mat' است که در نرم افزار MATLAB استفاده میشود . برای خواندن این نوع فایل در پایتون، از کتابخانه 'scipy.io' استفاده می کنیم. دستور زیر:

در اینجا کد به زبان پایتون آورده شده است:

```
\ data = scipy.io.loadmat("/content/252.mat")
\ print("Keys in MAT file:", data.keys())
```

خروجی:

```
Keys in MAT file: dict_keys(['__header__', '__version__', '__globals__',
'X252_DE_time', 'X252_FE_time', 'X252RPM'])
```

در اینجا، فایل MAT ۲۵۲ خوانده می شود و در متغیر data ذخیره می شود. با data.keys() می توانیم کلیدهای موجود در فایل را مشاهده کنیم.

پس اجزای آن شامل موارد زیر است:

header → فایل فرمت به مربوط اطلاعات

version → فایل نسخه ی MAT

globals → فایل در عمومی متغیرهای

X۲۵۲_DE_time → درایو انتهای لرزش سیگنال

X۲۵۲_FE_time → فن انتهای لرزش سیگنال

X۲۵۲RPM → چرخش سرعت (RPM)

بررسی نوع داده خوانده شده با دستور زیر:

```
\ print(type(data))
```

خروجی:

```
\ <class 'dict'>
```

نتیجه: داده ی خوانده شده از نوع دیکشنری (dict) است.

از بین سیگنال های موجود، می توانیم یکی از ستون های X۲۵۲_DE_time یا X۲۵۲_FE_time را انتخاب کنیم. در اینجا، سیگنال X۲۵۲_DE_time را انتخاب می کنیم. انتخاب سیگنال لرزش انتهای درایو و ذخیره در متغیر:

```
\ signal = data["X252_DE_time"].flatten()
```

ب) نمایش سیگنال :

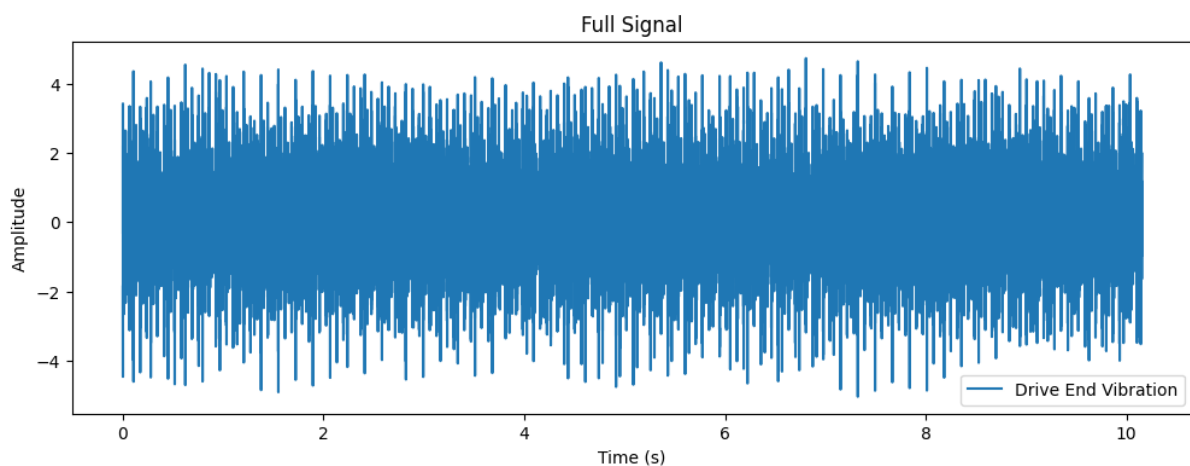
فرکانس نمونه برداری (Sampling Frequency): گفته شده که ۴۸KHz است. یعنی هر ثانیه ۴۸۰۰۰ نمونه ثبت شده است.

پس اگر N نمونه داشته باشیم، زمان معادل آن را می توان به این شکل محاسبه کرد:

$$\text{index fs time} = \frac{\text{index}}{fs} \quad \text{index} = \text{time}$$

کد برای رسم کل سیگنال:

```
۱ fs = 48000
۲ N = len(signal)
۳
۴ time = np.arange(N) / fs
۵
۶ plt.figure(figsize=(12, 4))
۷ plt.plot(time, signal, label="Drive End Vibration")
۸ plt.xlabel("Time (s)")
۹ plt.ylabel("Amplitude")
۱۰ plt.title("Full Signal")
۱۱ plt.legend()
۱۲ plt.show()
```

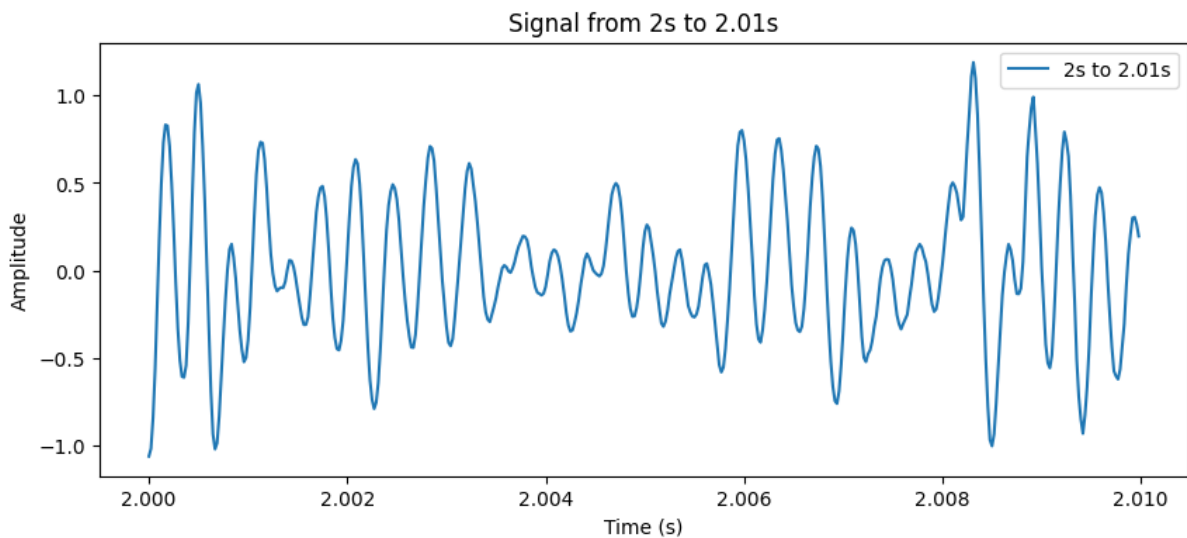


نمایش سیگنال در بازه‌ی $2s$ تا $2.01s$
محاسبه‌ی تعداد نمونه‌ها در این بازه: بازه‌ی $2s$ تا $2.01s$ یعنی 10 میلی‌ثانیه که معادل:

$$\text{Number of samples} = 0.01 \times 48000 = 480 \text{ samples}$$

کد برای رسم این بخش:

```
1 start_idx = int(2 * fs) #start
2 end_idx = start_idx + 480 #end sample
3
4 plt.figure(figsize=(10, 4))
5 plt.plot(time[start_idx:end_idx], signal[start_idx:end_idx], label="2s to
6         2.01s")
7 plt.xlabel("Time (s)")
8 plt.ylabel("Amplitude")
9 plt.title("Signal from 2s to 2.01s")
10 plt.legend()
11 plt.show()
```



ج) تحلیل فرکانسی :

تبدیل فوریه (FFT) برای پیدا کردن فرکانس‌های موجود در سیگنال انجام می‌شود.

مراحل:

تبدیل سیگنال به حوزه‌ی فرکانس با `np.fft.fft()`

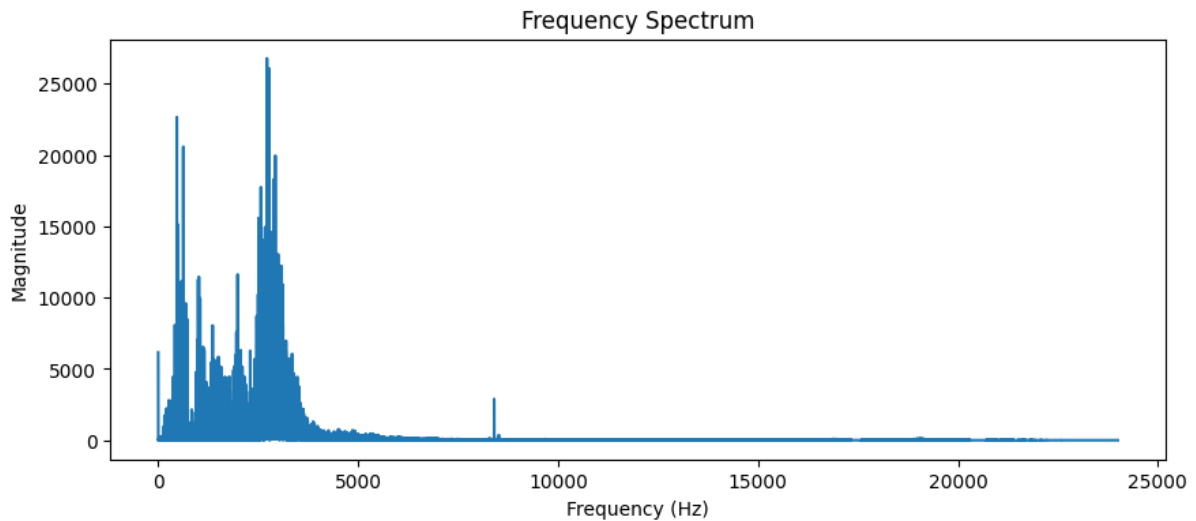
محور فرکانس را با `np.fft.fftfreq()` ایجاد کنیم

قدرتمندی طیف را با `np.abs()` محاسبه کنیم

کد:

```
۱ fft_vals = np.fft.fft(signal)
۲ freqs = np.fft.fftfreq(len(signal), d=1/fs)
۳
۴ magnitude = np.abs(fft_vals)
۵
۶ plt.figure(figsize=(10, 4))
۷ plt.plot(freqs[:N//2], magnitude[:N//2])
۸ plt.xlabel("Frequency (Hz)")
۹ plt.ylabel("Magnitude")
۱۰ plt.title("Frequency Spectrum")
۱۱ plt.show()
```

۱. رسم طیف فرکانسی

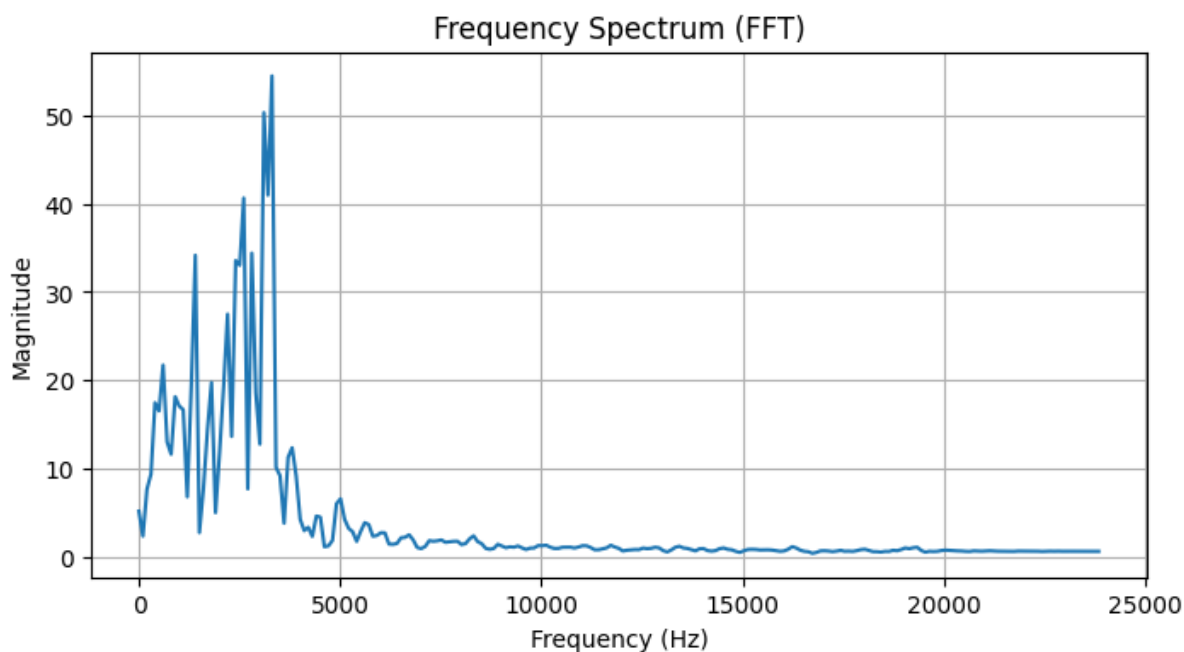


حال یک تابع می‌نویسم که: یک بازه زمانی دلخواه از سیگنال را به عنوان ورودی بگیرد، تبدیل فوریه (FFT) را روی آن اجرا کند و طیف فرکانسی را نمایش دهد.
 کد تابع برای نمایش طیف فرکانسی با تبدیل فوریه:

```

1 selected_signal = data['X252_DE_time'].squeeze()
2 def plot_frequency_spectrum(signal, sample_rate, start_time=0,
   end_time=None):
3
4     total_samples = len(signal)
5
6     if end_time is None:
7         end_time = total_samples / sample_rate
8
9     start_index = int(start_time * sample_rate)
10    end_index = int(end_time * sample_rate)
11
12    selected_signal = signal[start_index:end_index]
13
14    fft_result = np.fft.fft(selected_signal)
15    freq = np.fft.fftfreq(len(selected_signal), d=1/sample_rate)
16
17    positive_freqs = freq[:len(freq)//2]
18    magnitude = np.abs(fft_result[:len(freq)//2])
19
20    plt.figure(figsize=(8,4))
21    plt.plot(positive_freqs, magnitude)
22    plt.xlabel("Frequency (Hz)")
23    plt.ylabel("Magnitude")
24    plt.title("Frequency Spectrum (FFT)")
25    plt.grid()
26    plt.show()
27
28 plot_frequency_spectrum(selected_signal, sample_rate=48000, start_time=2,
   end_time=2.01)

```



`squeeze()` باعث می‌شود که اگر داده‌های ما به صورت ماتریس ذخیره شده باشند، به یک آرایه یک‌بعدی تبدیل شوند که پردازش آن آسان‌تر است. برای نمایش طیف فرکانسی سیگنال، یک تابع عمومی نوشته‌ایم که می‌تواند هر سیگنال را در یک بازه‌ی مشخص دریافت کرده و تبدیل فوریه‌ی آن را محاسبه کند.

`np.fft.fft()` تبدیل فوریه را روی سیگنال اعمال می‌کند.

`np.fft.fftfreq()` فرکانس‌های متناظر با خروجی تبدیل فوریه را محاسبه می‌کند.

فقط نیمه‌ی مثبت فرکانس‌ها نمایش داده می‌شود، زیرا در تبدیل فوریه، قسمت منفی فرکانس‌ها تقارن دارد و اطلاعات اضافه‌ای ندارد.

نمودار خروجی با محور افقی فرکانس (Hz) و محور عمودی دامنه (Magnitude) رسم می‌شود.

۲. پیدا کردن فرکانس غالب (بیشترین انرژی در طیف)

```
1 dominant_freq = freqs[np.argmax(magnitude[:N//2])]
2 print("Dominant Frequency:", dominant_freq, "Hz")
```

خروجی:

Dominant Frequency: 2719.5640398535857 Hz

د، ه (تقسیم بندی سیگنال

مراحل:

سیگنال را به قطعات ۱۲۸ تایی بشکنیم.

هر قطعه را در یک سطر از `numpy array` ذخیره کنیم.

مقدار هیپوتنوسی (Norm) هر قطعه را حساب کنیم.

```
1 import numpy as np
2 segment_size = 128
3 num_segments = len(signal) // segment_size
4
5 segments = signal[:num_segments * segment_size].reshape(num_segments,
6 segment_size)
7
8 segment_norms = np.linalg.norm(segments, axis=1)
9
10 print("Number of pieces:", num_segments)
11 print("Some sample hypotenuse values:", segment_norms[:5])
```

تعداد قطعات: ۳۸۰۷

چند مقدار هیپوتنوسی نمونه:

[18.75708897 9.5150878 12.12775435 12.15589371 8.60514374]

محاسبه‌ی مقدار هیپوتنوسی هر قطعه و ذخیره آن

(پاسخ به: "مقدار هیپوتنوسی قطعه‌ها را ذخیره کنید.")

```
1 segment_norms = np.linalg.norm(segments, axis=1)
2 print("First 10 segment norms:", segment_norms[:10])
```

خروجی:

First 10 segment norms: [18.75708897 9.5150878 12.12775435 12.15589371 8.60514374
6.47621605 3.01400884 4.63807857 3.21532278 1.87234841]

DataFrame در ذخیره قطعات

مراحل: هر قطعه را در یک `pandas.DataFrame` ذخیره کنیم. مقدار هیپوتنوسی هر قطعه را به آن اضافه کنیم.

```
1 import pandas as pd
2 df_segments = pd.DataFrame(segments)
3
4 df_segments["Norm"] = segment_norms
5
6 print(df_segments.head())
```

نمایش ۱۰ قطعه‌ی مضرب ۱۳ در یک نمودار

با استفاده از DataFrame ایجادشده، ۱۰ قطعه‌ی مضرب ۱۳ را در یک نمودار روی هم نمایش دهید. نمودارها باید رنگ‌های مختلف داشته باشند و نمودار نهایی نیز برجسب‌دار باشد. (تمامی این مراحل باید در یک حلقه for انجام شود).

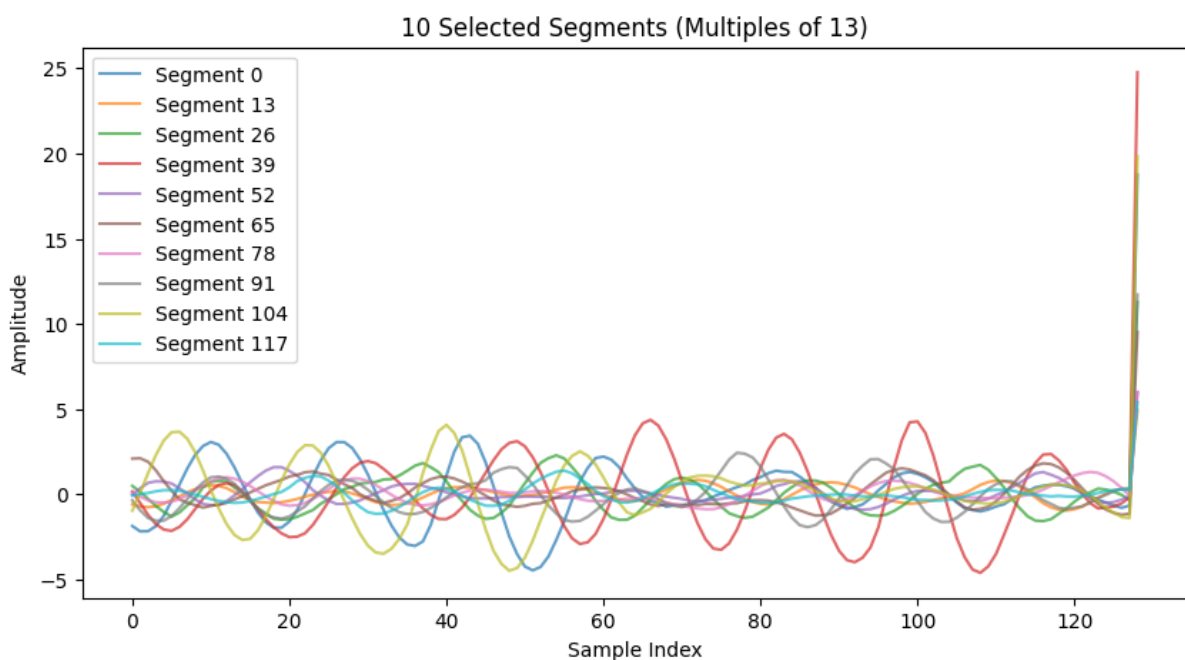
توضیح:

۱۰ قطعه‌ی اولی که ایندکس آنها مضرب ۱۳ است را پیدا می‌کنیم. هرکدام را روی نمودار رسم می‌کنیم و از رنگ‌های متفاوت استفاده می‌کنیم.

```

۱ import matplotlib.pyplot as plt
۲ selected_segments = df_segments.iloc[:,13][:10].values
۳
۴ plt.figure(figsize=(10, 5))
۵ for i, segment in enumerate(selected_segments):
۶     plt.plot(segment, label=f"Segment {i*13}", alpha=0.7)
۷
۸ plt.xlabel("Sample Index")
۹ plt.ylabel("Amplitude")
۱۰ plt.title("10 Selected Segments (Multiples of 13)")
۱۱ plt.legend()
۱۲ plt.show()

```



(استخراج ویژگی (Feature Extraction)

```

۱ features = [extract_features(segment) for segment in segments]
۲
۳ df_features = pd.DataFrame(features)
۴
۵ print(df_features.head())
۶ df_features.to_csv("extracted_features.csv", index=False)
۷ print("File saved: extracted_features.csv")
۸ from google.colab import files
۹
۱۰ files.download("extracted_features.csv")

```

- تقسیم سیگنال به قطعات ۱۲۸ تایی و ذخیره در numpy
- ذخیره‌ی قطعات در pandas.DataFrame
- رسم ۱۰ قطعه‌ی مضرب ۱۳ روی یک نمودار
- نوشتن تابعی برای استخراج ویژگی‌های آماری
- ذخیره‌ی ویژگی‌های محاسبه‌شده در DataFrame و خروجی CSV

Iris Dataset II

مقدمه

تحقیق درباره دیتاست معروف Iris

دیتاست Iris که اولین بار توسط Sir R.A. Fisher استفاده شد، از مقاله‌ی فیشر گرفته شده است. این مجموعه داده همان مجموعه‌ای است که در R استفاده شده، اما کمی با نسخه‌ای که در UCI Machine Learning Repository وجود دارد، متفاوت است؛ زیرا نسخه‌ی UCI دو نقطه داده‌ی اشتباه دارد. این مجموعه داده یکی از شناخته‌شده‌ترین دیتاست‌ها در ادبیات شناسایی الگو (Pattern Recognition) است. مقاله‌ی فیشر یک مقاله‌ی کلاسیک در این حوزه محسوب می‌شود و هنوز هم به آن ارجاع داده می‌شود (برای مثال، در کتاب (Duda & Hart)).

در این گزارش، به بررسی Iris Dataset پرداخته می‌شود که یکی از معروف‌ترین مجموعه‌داده‌ها در زمینه یادگیری ماشین و تحلیل داده است. این مجموعه داده شامل ویژگی‌های گل‌های Iris در سه گونه‌ی مختلف است و برای طبقه‌بندی استفاده می‌شود. هدف اصلی این بررسی، درک توزیع داده‌ها، تحلیل ویژگی‌ها و بررسی امکان‌پذیری طبقه‌بندی آنها بر اساس ویژگی‌هایشان است.

شرح داده‌ها

مجموعه داده‌ی Iris شامل ۱۵۰ نمونه از سه گونه‌ی گل Iris است که به طور مساوی (۵۰ نمونه برای هر گونه) توزیع شده‌اند. هر نمونه دارای ۴ ویژگی عددی است:

- طول کاسبرگ (Sepal Length)
- عرض کاسبرگ (Sepal Width)
- طول گلبرگ (Petal Length)
- عرض گلبرگ (Petal Width)

علاوه بر این ویژگی‌ها، یک برچسب دسته‌ای (Class Label) وجود دارد که نشان می‌دهد هر نمونه متعلق به کدام گونه از گل‌های Iris است.

- Setosa
- Versicolor
- Virginica

تحلیل آماری و بررسی ویژگی‌ها

توزیع داده‌ها

با بررسی داده‌ها، مشخص می‌شود که ویژگی‌های مختلف گل‌ها دارای **توزیع‌های متفاوتی** هستند. برخی ویژگی‌ها مانند **طول گلبرگ (Petal Length)** در گونه‌های مختلف تفاوت‌های واضحی دارند، در حالی که برخی ویژگی‌های دیگر مانند **عرض کاسبرگ (Sepal Width)** ممکن است هم‌پوشانی بیشتری بین گونه‌ها داشته باشند.

نمایش گرافیکی داده‌ها

برای درک بهتر داده‌ها، می‌توان از نمودارهای **توزیع**، **هیستوگرام**، **جعبه‌ای (Boxplot)** و **scatter plot** استفاده کرد. به عنوان مثال:

- نمودار پراکندگی (scatter plot) نشان می‌دهد که داده‌های گونه‌ی **Setosa** به خوبی از دو گونه‌ی دیگر جدا شده‌اند.
- نمودار هیستوگرام به ما کمک می‌کند توزیع هر ویژگی را به صورت جداگانه ببینیم.
- نمودار جعبه‌ای (Boxplot) میزان پراکندگی و نقاط پرت (Outliers) را نشان می‌دهد.

تحلیل تفکیک‌پذیری کلاس‌ها

یکی از مهم‌ترین سوالات در این تحلیل این است که آیا می‌توان گونه‌های مختلف را بر اساس ویژگی‌هایشان به درستی از هم تفکیک کرد؟

- **Setosa** به وضوح از بقیه گونه‌ها متمایز است.
- **Versicolor** و **Virginica** هم‌پوشانی بیشتری دارند، ولی در برخی ویژگی‌ها مانند طول و عرض گلبرگ تفاوت مشخصی دارند.

نتیجه‌گیری

مجموعه داده‌ی **Iris** یک مثال عالی برای تحلیل داده و یادگیری ماشین است. بررسی ویژگی‌های این داده‌ها نشان می‌دهد که گونه‌های مختلف گل‌ها را می‌توان با تحلیل ویژگی‌هایشان از هم تفکیک کرد. این مجموعه داده به دلیل ساختار ساده ولی کاربردی‌اش، یکی از اولین گزینه‌ها برای یادگیری مفاهیم داده‌کاوی، طبقه‌بندی و تحلیل آماری محسوب می‌شود.

خلاصه‌ی کلی گزارش این بخش

- معرفی داده‌ها و ویژگی‌های آنها
- تحلیل توزیع داده‌ها و نمایش بصری
- بررسی امکان تفکیک‌پذیری کلاس‌ها
- نتیجه‌گیری درباره‌ی استفاده‌ی این داده‌ها در یادگیری ماشین

```

1 # Import necessary libraries
2 from sklearn import datasets # Load built-in datasets from scikit-learn
3 import pandas as pd # Data manipulation and analysis
4 import numpy as np # Numerical computations
5 from sklearn.model_selection import train_test_split # Splitting dataset
   into train and test
6 import matplotlib.pyplot as plt # Data visualization
7 import seaborn as sns # Advanced visualization
8 from mpl_toolkits.mplot3d import Axes3D # 3D plotting
9
10 # Load the Iris dataset
11 iris = datasets.load_iris()
12
13 # Print dataset description
14 print(iris.DESCR)
15
16 # Convert dataset to a Pandas DataFrame with feature names as columns
17 data = pd.DataFrame(data=iris.data, columns=iris.feature_names)
18 # Add the target variable (species label) to the DataFrame
19 data['target'] = iris.target
20
21 # Split the dataset into training (70%) and testing (30%) sets
22 data_train, data_test = train_test_split(data, test_size=0.3,
   random_state=42)
23
24 # Label the sets to differentiate between train and test data
25 data_train['set'] = 'train'
26 data_test['set'] = 'test'
27
28 # Combine training and test datasets into a single DataFrame
29 data_combined = pd.concat([data_train, data_test])
30
31 # Display the first five rows of the combined dataset
32 print(data_combined.head())
33
34 # 2D Scatter plot using the first two features
35 plt.figure(figsize=(8, 6))
36 sns.scatterplot(
37     x=data_combined[iris.feature_names[0]], # Sepal length
38     y=data_combined[iris.feature_names[1]], # Sepal width
39     hue=data_combined['target'], # Color by species
40     palette='viridis' # Color theme
41 )
42 plt.xlabel(iris.feature_names[0]) # Label x-axis
43 plt.ylabel(iris.feature_names[1]) # Label y-axis
44 plt.title("Scatter of samples based on two characteristics") # Title of
   the plot
45 plt.legend(iris.target_names) # Add species labels as legend
46 plt.show()
47
48 # 3D Scatter plot using the first three features
49 fig = plt.figure(figsize=(10, 8))
50 ax = fig.add_subplot(111, projection='3d')
51
52 # Loop through each species and plot them in 3D
53 for i, species in enumerate(iris.target_names):
54     subset = data_combined[data_combined['target'] == i]
55     ax.scatter(subset[iris.feature_names[0]],

```

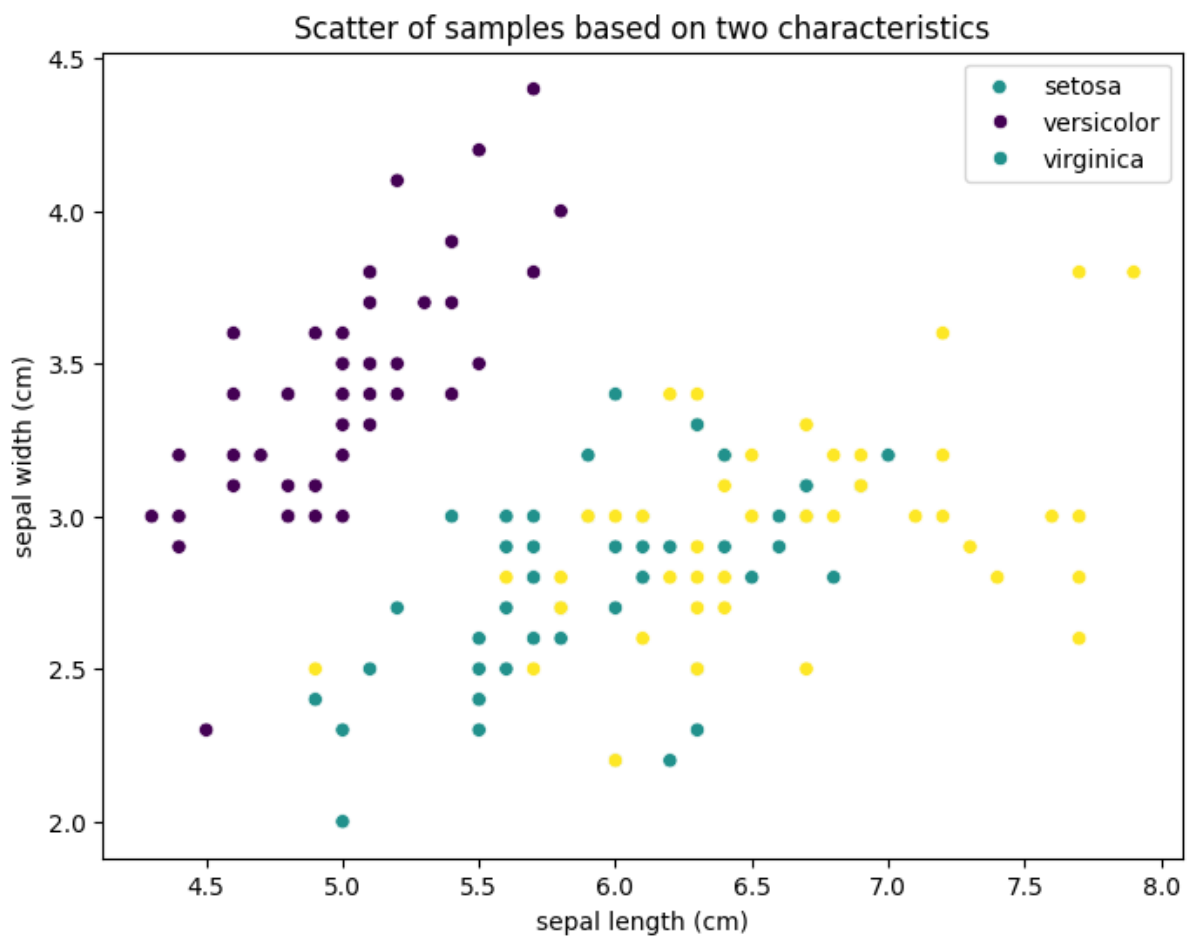
```

subset[iris.feature_names[1]], subset[iris.feature_names[2]],
label=species)

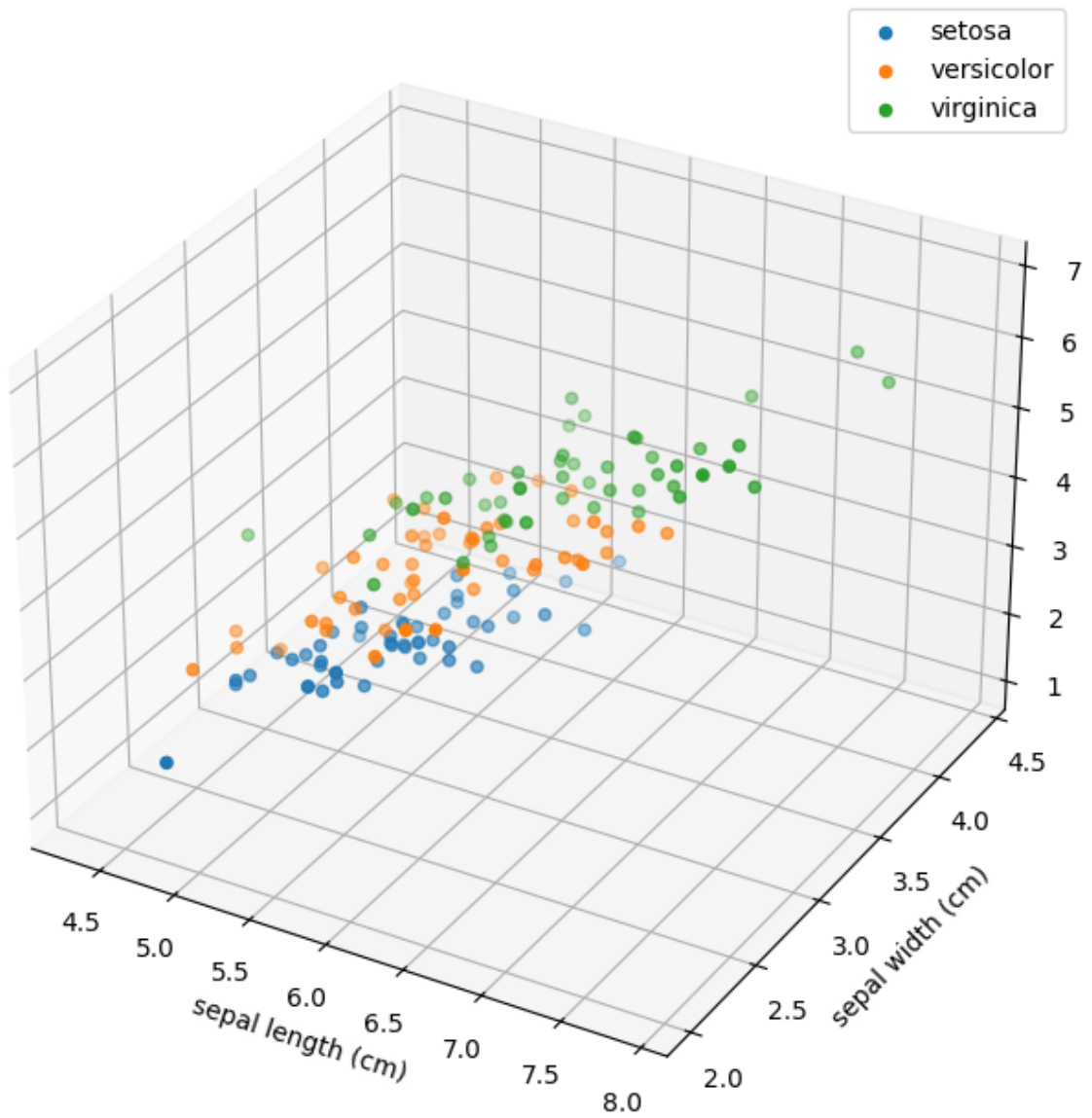
۵۶
۵۷ # Label axes
۵۸ ax.set_xlabel(iris.feature_names[0])
۵۹ ax.set_ylabel(iris.feature_names[1])
۶۰ ax.set_zlabel(iris.feature_names[2])
۶۱
۶۲ # Title and legend
۶۳ ax.set_title("Three-dimensional data dispersion")
۶۴ ax.legend()
۶۵ plt.show()
۶۶
۶۷ # Heatmap of feature correlations
۶۸ plt.figure(figsize=(10, 8))
۶۹ sns.heatmap(data_combined[iris.feature_names].corr(), annot=True,
cmap='coolwarm', linewidths=0.5)
۷۰ plt.title("Feature correlation heat map")
۷۱ plt.show()
۷۲
۷۳ # Kernel Density Estimation (KDE) plot for each feature
۷۴ for feature in iris.feature_names:
۷۵     plt.figure(figsize=(8, 6))
۷۶     sns.kdeplot(data_train[feature], label='Train', fill=True, alpha=0.5)
        # Train data distribution
۷۷     sns.kdeplot(data_test[feature], label='Test', fill=True, alpha=0.5) #
        Test data distribution
۷۸     plt.xlabel(feature)
۷۹     plt.title(f"Feature probability density {feature}") # Title
        indicating which feature is displayed
۸۰     plt.legend()
۸۱     plt.show()
۸۲
۸۳ # Discretization of the first feature (sepal length) into categorical bins
۸۴ bins = [data_combined[iris.feature_names[0]].min(), 5.5, 6.5,
        data_combined[iris.feature_names[0]].max()] # Define bin edges
۸۵ labels = ['short', 'medium', 'tall'] # Labels for each category
۸۶
۸۷ # Create a new categorical column based on the feature value
۸۸ data_combined['discrete_feature'] =
        pd.cut(data_combined[iris.feature_names[0]], bins=bins, labels=labels)
۸۹
۹۰ # Display first five rows with the new categorical feature
۹۱ print(data_combined[['discrete_feature', iris.feature_names[0]]].head())
۹۲
۹۳ # Extract only the data for the "Setosa" species (target = 0)
۹۴ setosa_data = data_combined[data_combined['target'] == 0]
۹۵
۹۶ # Display summary statistics for the Setosa subset
۹۷ print(setosa_data.describe())

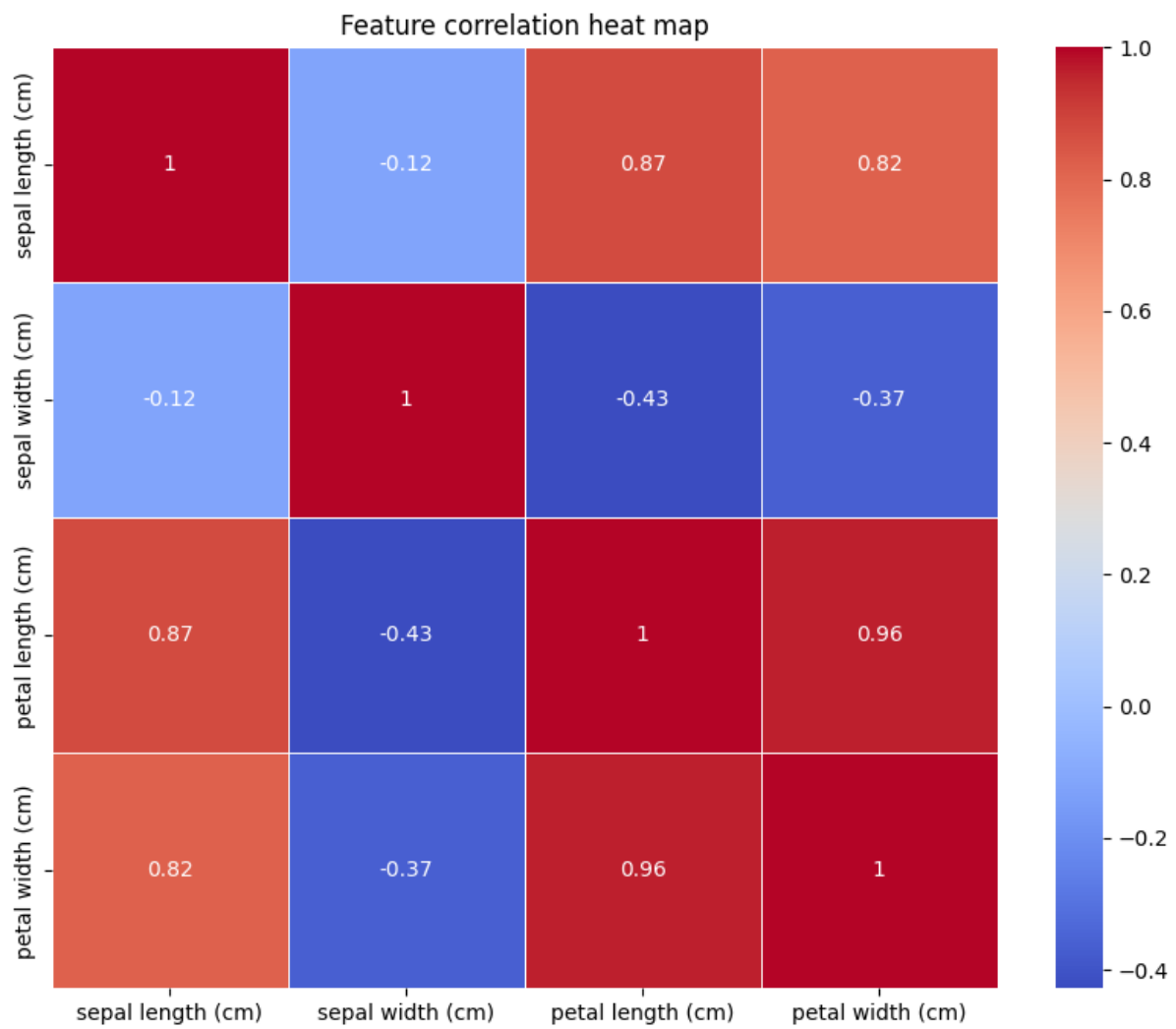
```

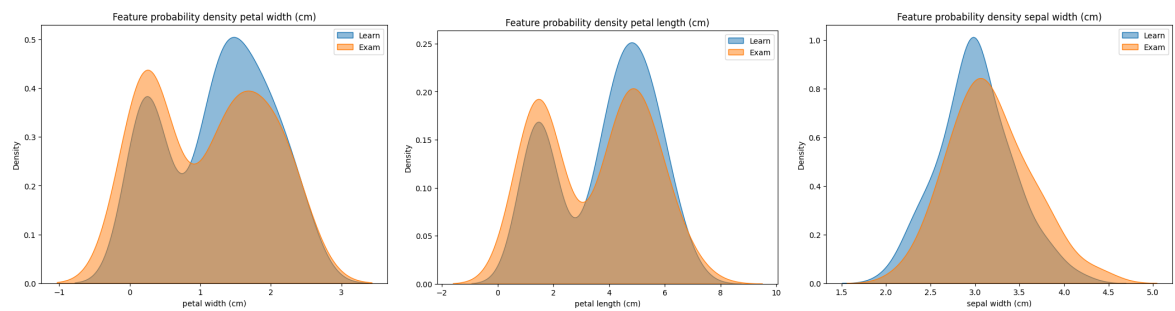
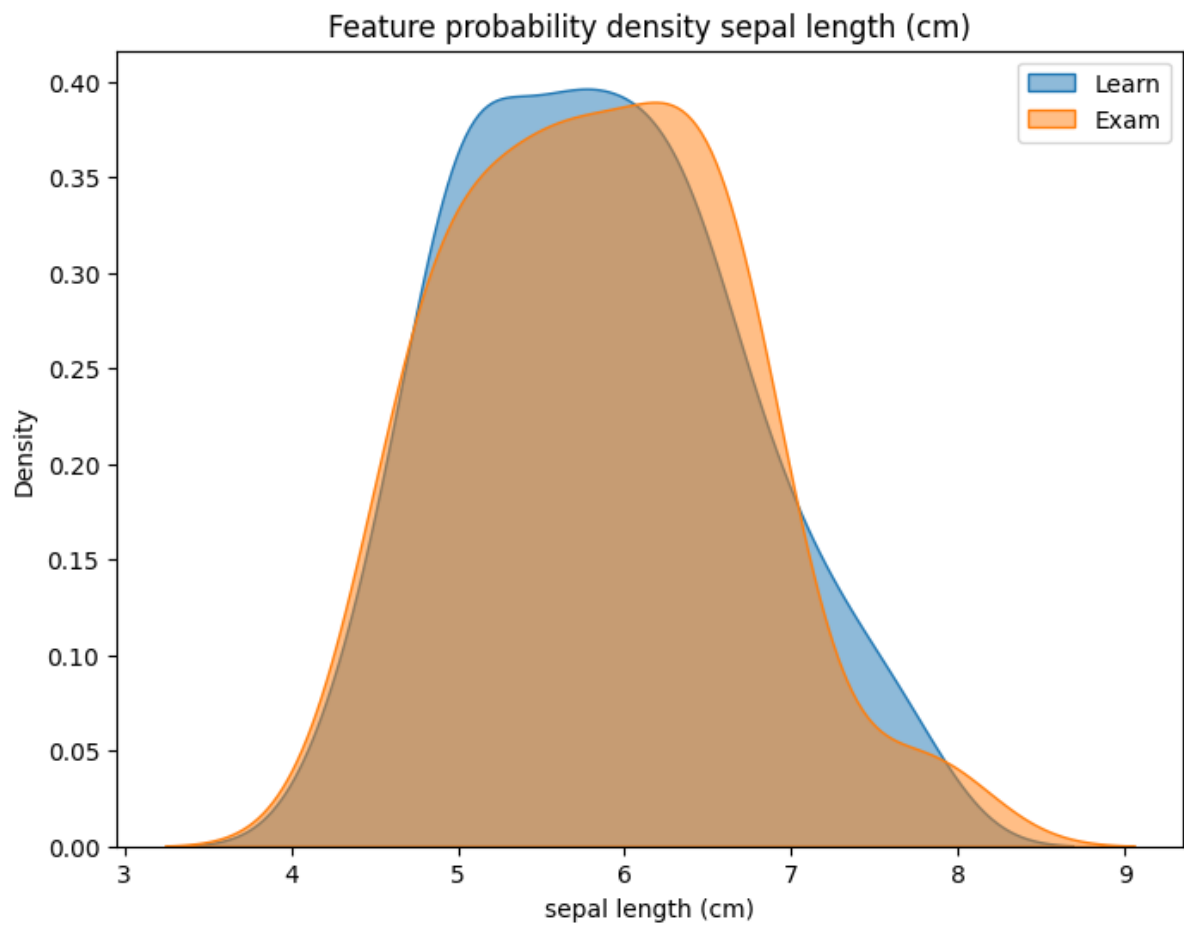
ممنون از توجه شما.



Three-dimensional data dispersion







(ج) تصویر ۳

(ب) تصویر ۲

(آ) تصویر ۱