

دانشگاه صنعتی امیر کبیر
(پلی تکنیک تهران)

درس : کلان داده

نام و نام خانوادگی : فاطمه توکلی

۴۰۰۱۳۱۰۱۶

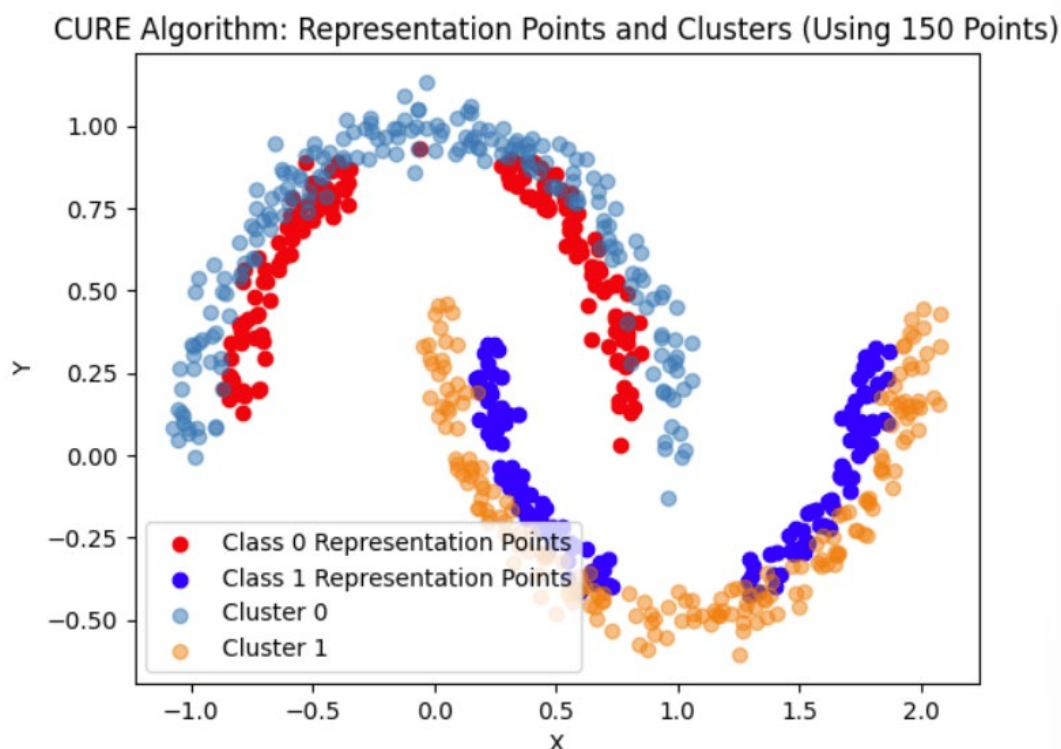
تمرین شماره ۰۲

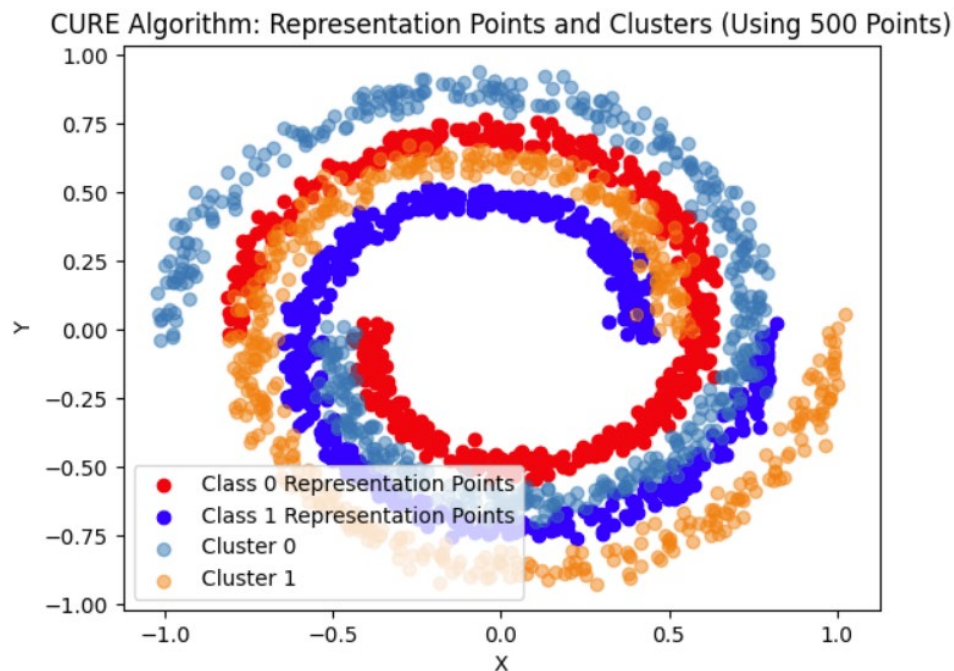
الف) در ابتدا به صورت رندوم از بین نقاط تولید شده تعدادی را انتخاب میکنیم و به عنوان نمونه با این تعداد داده پیش میرویم و سپس بر روی آن ها با استفاده از **AgglomerativeClustering** دسته بندی سلسله مراتبی روی داده ها انجام داده، سپس برای انتخاب نقاط **representative** در هر دسته به صورت رندوم یک نقطه شروع انتخاب میکنیم و در صورتی که تعداد نقاط **representative** به حد دلخواه که ورودی تابع است نرسیده باشد به صورت تکرار شونده نقطه بعدی به صورتی انتخاب میشود که بیشترین فاصله اقلیدسی با نقطه فعلی را داشته باشد. و در نهایت تمامی نقاط ۲۰ درصد به سمت مرکز هر دسته شیفت داده میشوند.

در مجموعه داده اول برای اینکه دسته بندی به درستی انجام شود از ۱۰۰۰ داده تولید شده ۴۰۰ داده به عنوان نمونه انتخاب شده و تعداد نماینده ها ۱۵۰ نقطه است

در مجموعه داده دوم از ۲۰۰۰ داده تولید شده ۱۰۰۰ به عنوان نمونه انتخاب شده و تعداد نماینده ها ۵۰۰ نقطه میباشد.

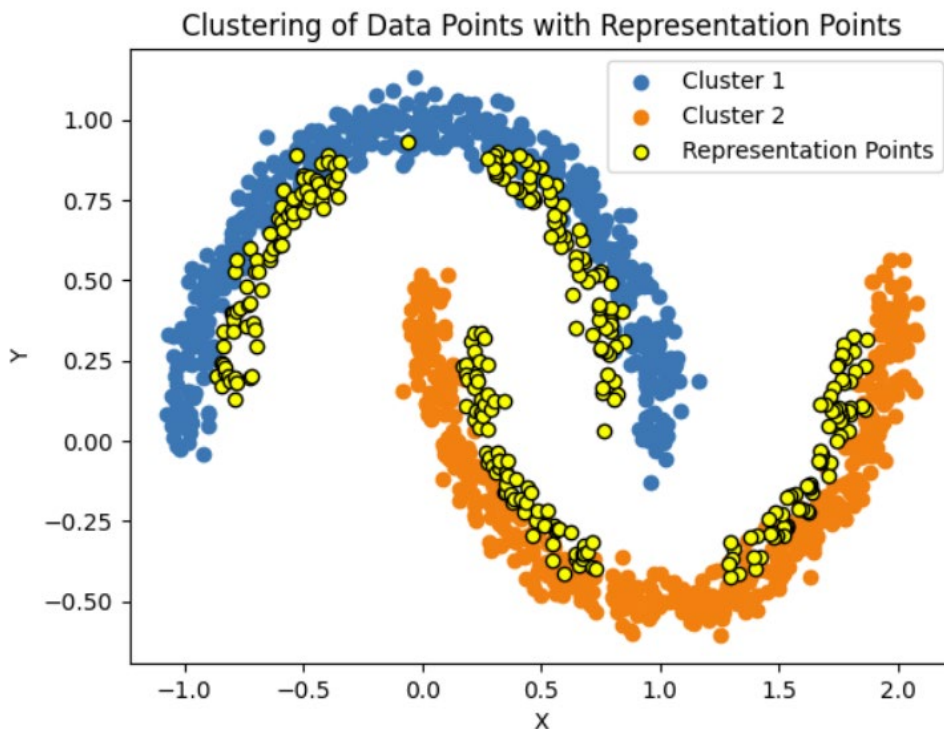
ب) نتیجه برای دو مجموعه داده آورده شده به صورت زیر است:

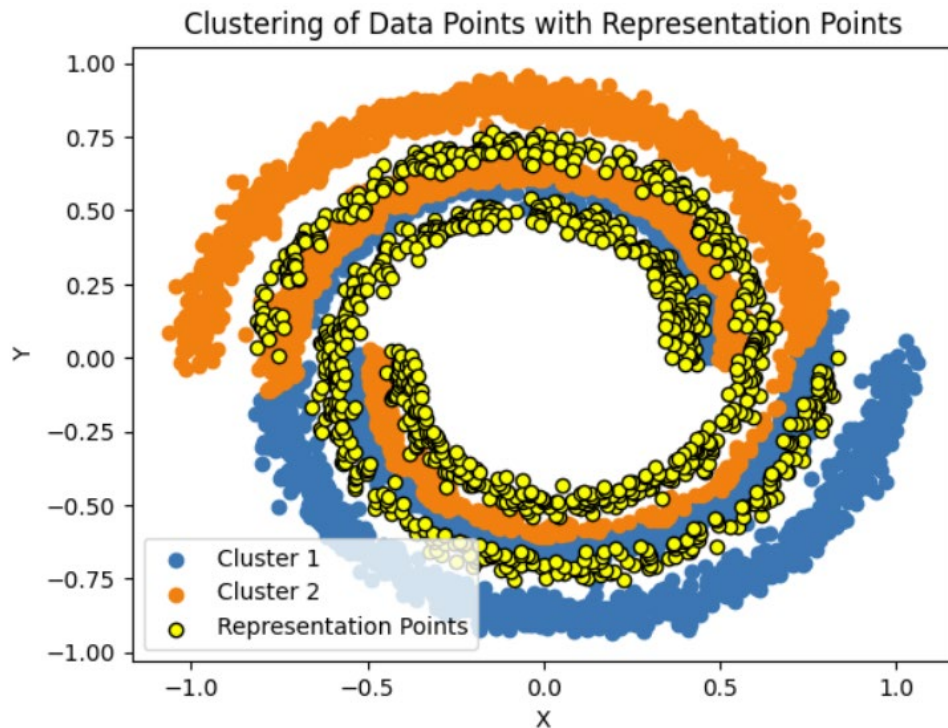




ج) برای دسته بندی نقاط همانطور که در جزوه نیز آورده شده نزدیک نقطه نماینده به هر ورودی را یافته و ورودی به آن دسته اضافه میشود

نتایج برای دو مجموعه داده به صورت زیر است و همانطور که مشاهده میشود برای مجموعه داده spiral به خوبی جواب نداده است





(د)

الگوریتم BFR (Boundary-based Fast Representative) یک الگوریتم خوشه‌بندی که به صورت مرحله‌ای عمل می‌کند و در هر مرحله تلاش می‌کند نماینده‌های خوشه‌ها را با استفاده از مرزهای خوشه‌ها بهبود دهد. در ابتدا k دسته در نظر می‌گیرد و داده را به صورت دسته دسته از حافظه می‌خواند و هر داده جدید را به یکی از دسته‌ها نسبت می‌دهد اگر که فاصله آن‌ها از یک مقدار آستانه کوچکتر باشد. و اگر فاصله از مقدار آستانه کوچکتر نباشد آن را در دسته دیگر به نام **compressed sets** قرار می‌دهد و در مرحله بعدی تلاش بر ادغام این **compress sets** ها از مرحله قبلی دارد و این کار را تکرار می‌کند تا زمانی که تمام نقاط دسته بندی شوند. کار این الگوریتم خلاصه کردن داده‌های حجیم است و مسئله انتخاب معیار فاصله و تصمیم‌گیری راجع به ادغام مجموعه‌ها بر نتایج اثر دارد

عیب این روش این است که واریانس داده‌ها را حول محور X یا Y محدود می‌کند و داده باید در جهت یکی از محورها توزیع شده باشد.

مزایای الگوریتم BFR عبارتند از:

۱. کارایی BFR برای پردازش مجموعه داده‌های بزرگ بهینه شده است و می‌تواند با سرعت بالا و موازی اجرا شود.
۲. کنترل بر روی ریز خوشه‌ها BFR می‌تواند به طور هوشمند ریز خوشه‌ها را تعیین کند و در صورت لزوم تعداد آنها را کاهش دهد تا از بهبود کیفیت خوشه اطمینان حاصل شود.
۳. الگوریتم BFR می‌تواند با خوشه‌بندی مجدد خوشه‌ها، با استفاده از سرعت و کارایی بالای خود، از داده‌های قبلی زمانی که تغییرات رخ می‌دهد یا داده‌های جدید اضافه می‌شود، استفاده مجدد کند.

مزایای الگوریتم CURE عبارتند از:

۱. CURE در برابر نقاط پرت مقاوم است.

۲. با استفاده از نقاط نماینده برای تعریف خوشه ها، نتایج قابل تفسیر را به راحتی ارائه می دهد.

معایب الگوریتم BFR عبارتند از:

۱. BFR به شدت به تخمین دقیق مرزهای خوشه ای متکی است، که می تواند در مجموعه داده های پیچیده با خوشه های همپوشانی یا شکل نامنظم چالش برانگیز باشد.

۲. حساسیت به تنظیمات پارامتر: عملکرد BFR می تواند به انتخاب پارامترها مانند تعداد ریز خوشه ها حساس باشد.

معایب الگوریتم CURE عبارتند از:

۱. پیچیدگی محاسباتی: CURE به هزینه محاسباتی بالاتری در مقایسه با BFR نیاز دارد، به خصوص زمانی که با مجموعه داده های بزرگ سروکار داریم، به دلیل نیاز به محاسبه فاصله بین تمام نقاط داده و نمایندگان.

۲. وابستگی به انتخاب نماینده: کیفیت نتایج خوشه بندی در CURE به شدت به انتخاب نقاط نماینده بستگی دارد که می تواند ذهنی باشد و بر عملکرد کلی تأثیر بگذارد.

به طور خلاصه، BFR کارایی و کنترل را بر روی خوشه های کوچک ارائه می کند، در حالی که CURE استحکام و قابلیت تفسیر را فراهم می کند. با این حال، BFR به تخمین مرز و تنظیمات پارامتر حساس است، در حالی که CURE پیچیدگی محاسباتی بالاتری دارد و به انتخاب نماینده وابسته است. انتخاب بین دو الگوریتم به ویژگی های خاص مجموعه داده و اهداف مورد نظر کار خوشه بندی بستگی دارد.

ه) همانطور که دیدیم در مجموعه داده دوم نتوانست اینکار را به درستی انجام دهد یکی از دلایل آن میتواند ۲۰ درصد انتقال به سمت مرکز باشد که باعث شده نقاط نماینده درهم تنیده شوند و معیار فاصله به درستی نتواند دسته ها را تفکیک کند.

ی) روش elbow تکنیکی است که برای تعیین تعداد بهینه خوشه ها (k) در الگوریتم خوشه بندی k -means استفاده می شود. شامل ترسیم تعداد خوشه ها در برابر مجموع مربعات خطاها (SSE) است.

برای اعمال روش elbow، با انجام خوشه بندی k -means با طیفی از مقادیر k ، معمولاً از یک مقدار کوچک تا یک مقدار نسبتاً بزرگ شروع می شود. برای هر مقدار k ، SSE محاسبه شده و بر روی یک نمودار خطی رسم می شود. نمودار معمولاً روند کاهشی در SSE را با افزایش k نشان می دهد، زیرا خوشه های بیشتر می توانند بهتر با نقاط داده تناسب داشته باشند. با این حال، بعد از یک نقطه خاص، کاهش SSE شروع به صاف شدن می کند و شکل خم یا آرنج را در نمودار ایجاد می کند. تعداد بهینه خوشه ها اغلب به عنوان مقدار k در نقطه elbow در نظر گرفته می شود.

با توجه به استفاده از روش elbow برای تعیین k در الگوریتم K-NN، مناسب نیست زیرا الگوریتم K-NN یک الگوریتم طبقه بندی است، نه یک الگوریتم خوشه بندی. الگوریتم K-NN برچسب ها را به نقاط داده جدید بر اساس برچسب های k نزدیکترین همسایه آنها در داده های آموزشی اختصاص می دهد. این شامل خوشه بندی یا تعیین تعداد خوشه ها نیست.

در K-NN، مقدار k نشان دهنده تعداد همسایه هایی است که برای طبقه بندی در نظر گرفته می شوند، و معمولاً بر اساس دانش قبلی، اعتبارسنجی متقابل یا تخصص دامن انتخاب می شود. ربطی به مفهوم خوشه بندی یا روش آرنجی ندارد. بنابراین، روش elbow برای تعیین مقدار k در الگوریتم K-NN قابل اجرا نیست.

الف) یکی از دلایلی که این الگوریتم را برای این مسئله مناسب میکند این است که تضمین میکند false negative ندارد در حالیکه از حافظه کمی مصرف می کند و همچنین false positive آن توسط تعداد توابع هش مستقل و تعداد m, n قابل کنترل است و احتمال برخورد را میتوان صفر کرد.

ب) برای به دست آوردن اندازه جدول هش میتوان از فرمول زیر که در جزوه آمده استفاده کرد:

$$\text{false positive probability} = (1 - e^{-km/n})^k$$

که با جایگذاری میتوان به این نتیجه رسید که ۵ برابر جدول اصلی تعریف میشود.

```
# Example usage
primes = [17, 31, 47, 61] # List of prime numbers

bloom_filter = BloomFilter(1000000, 5000000, len(primes), hash_func_1)

# Load the dataset and add items to the Bloom filter using the specified hash function
dataset_path = "user_dataset.csv" # Replace with the actual path to your dataset file

with open(dataset_path, "r") as file:
    csv_reader = csv.reader(file)
    for row in csv_reader:
        username = row[0] # Assuming the username is in the first column
        bloom_filter.add(username)

# Get the bit array
bit_array = bloom_filter.get_bit_array()
```

ابتدا مجموعه داده خوانده شده و لیست اعداد اول تعریف شده به عنوان ورودی m, n تعداد تابع هش، و همچنین نوع تابع هش را میگیرد و تک تک username ها را خوانده و خانه های جدول هش را با توجه به ۴ تابع هش مقداردهی میکند.

```
def hash_func_1(item, i):
    c = [ord(ch) for ch in item]
    # Select prime number from the list based on the hash function index
    p = primes[i % len(primes)]
    l = len(c)
    min_c = min(c)
    max_c = max(c)
    index = (min_c + math.prod(c)) * p + ((sum(c)) * pow(p, (1 // 2))) + (max_c * pow(p, (1 - 1)))
    return index
```

تابع هش اول به این صورت تعریف می شود که ابتدا کاراکترهای موجود در رشته را به ارقام تبدیل میکند و سپس از لیست اعداد اول داده شده به ترتیب تابع هش را محاسبه میکند و خروجی را برمیگرداند.

```
index = self.hash_func(item, i) % self.m # Get the index for each hash function
self.bit_array[index] = True # Set the corresponding bit to 1
```

خروجی نظیر را بر m باقیمانده میگیرد و در آرایه مقدار آن را برابر با True در نظر میگیرد.

ج) تنها تفاوت در پیاده سازی تابع هش خواهد بود و بقیه موارد مانده قبل می باشد:

```
def hash_func_2(item, i):
    c = [ord(ch) for ch in item]
    p = primes[i % len(primes)]
    l = 0
    index = 0
    for x in c :
        index = index + (x * pow(p,l))
        l = l+1
    return index
```

د) برای اینکار username ها را به ترتیب میخوانیم و سپس خروجی تک تک توابع هش را بررسی میکنم اگر تمامی آن ها True بود یعنی آزاد نیست:

```
def test(self, item):
    hash_results = []
    for i in range(self.k):
        index = self.hash_func(item, i) % self.m # Get the index for each hash function
        hash_results.append(self.bit_array[index])
    if all(hash_results):
        return True
    else:
        print('{} is not in the dataset'.format(item))
        return False
```

و بدین صورت هر کدام در دیتاست نباشد در خروجی می آورد :

```
icolors is not in the dataset
eman88 is not in the dataset
PChrowsaway is not in the dataset
Anonymous132455 is not in the dataset
koi18 is not in the dataset
netskills001 is not in the dataset
vdubfast01 is not in the dataset
TheSnekMan is not in the dataset
askingmaggie is not in the dataset
```


مقدار false positive ,accuracy به صورت زیر است:

```
true positive rate is : 1.0  
Accuracy is : 0.87722
```

(۵)

```
koi18 is not in the dataset  
netskills001 is not in the dataset  
mn766132 is not in the dataset  
confusedouttamymind is not in the dataset  
vdubfast01 is not in the dataset  
TheSnekMan is not in the dataset  
askingmaggie is not in the dataset  
DankFlowers is not in the dataset  
WatchWhatHappensLive is not in the dataset  
Gamerulf is not in the dataset  
deppaotoko is not in the dataset
```

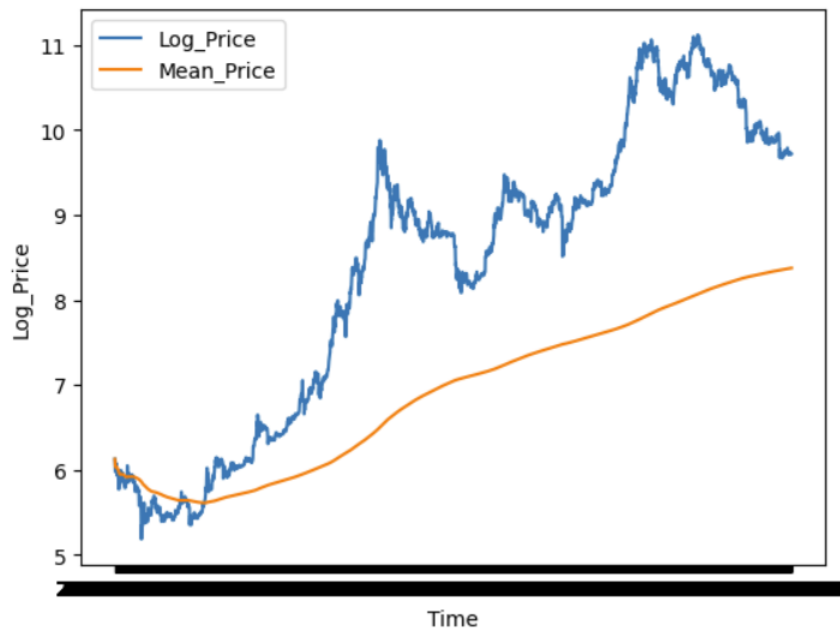
```
true positive rate is : 1.0  
Accuracy is : 0.91256
```

ی) خیر این مقدار بهتر است و ممکن است دلایل متفاوتی داشته باشد :

- ۱- مجموعه داده مورد استفاده ممکن است ویژگی های خاصی داشته باشد که دستیابی به نرخ مثبت کاذب مورد انتظار را برای فیلتر بلوم چالش برانگیز می کند.
- ۲- توابع هش انتخاب شده ممکن است برای مجموعه داده بهینه نباشد، که منجر به مثبت کاذب بالاتر می شود.
- ۳- اندازه آرایه بیت یا تعداد توابع هش استفاده شده در فیلتر Bloom ممکن است برای مجموعه داده مناسب نباشد و بر نرخ مثبت کاذب تأثیر بگذارد.
- ۴- تصادفی بودن مجموعه داده و انتخاب پرس و جو نیز می تواند به تغییرات در نرخ مثبت کاذب کمک کند.

الف) بعد از جدا کردن ستون close و همچنین لگاریتم گرفتن از آن ستون، ردیف به ردیف به صورت زیر میانگین را حساب کرده:

```
# Iterate through the log price row by row
for i in range(len(log_close)):
    mean = (mean * i + log_close[i]) / (i + 1)
    mean_values.append(mean)
```

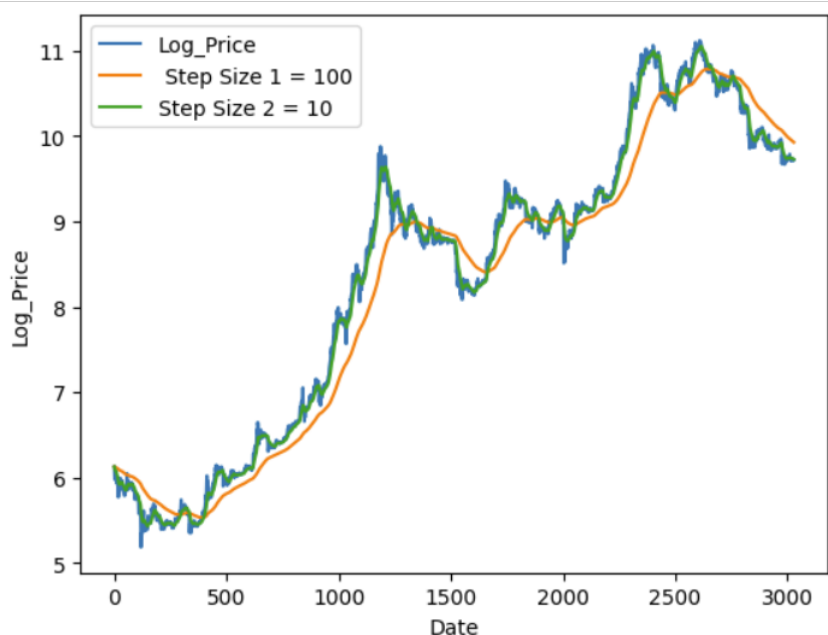


ب) دو مقدار برای step size و تابع لگاریتمی برای آن تعریف میکنیم:

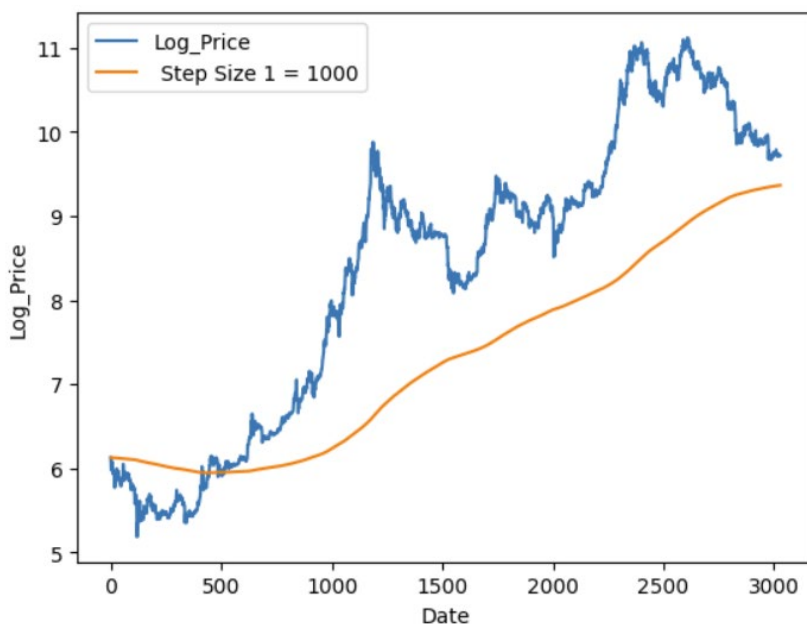
```
stepsize1 = 100
stepsize2 = 10
alpha1 = 1 - math.exp(-1 / stepsize1)
alpha2 = 1 - math.exp(-1 / stepsize2)
```

میانگین وزن دار به صورت زیر تعریف میشود:

```
for i in range(1, len(log_close)):
    new_price = log_close.iloc[i]
    ew_average1 = new_price * alpha1 + (1 - alpha1) * ew_average1
    ew_average_list1.append(ew_average1)
```



اگر بخواهیم شبیه قسمت اول شود باید $\text{stepsize} = 1000$ قرار داده شود:



ج) در این قسمت هم از میانگین وزن دار برای محاسبه اختلاف هر داده از واریانس استفاده کردیم هم برای میانگین گیری واریانس ها به صورت وزن دار عمل کردیم، و مقدار آستانه را به این صورت قرار داده که اگر ۴۰ بار به صورت متوالی قدر مطلق اختلاف واریانس از میانگین وزن دار واریانس بیشتر از ۰.۰۱ شد آن نقطه را به عنوان هشدار در نظر گرفتیم:

```
if abs(variance - var_mean) > 0.01:
    count += 1
    if count >= 40:
        points.append((i, new_price))
else:
    count = 0
```

و حاصل به صورت زیر است:

