

به نام خدا
دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

گزارش تمرین چهارم درس یادگیری ماشین

استاد درس: دکتر احسان ناظر فرد

دانشجو: فاطمه غلامزاده

۹۹۱۳۱۰۰۳

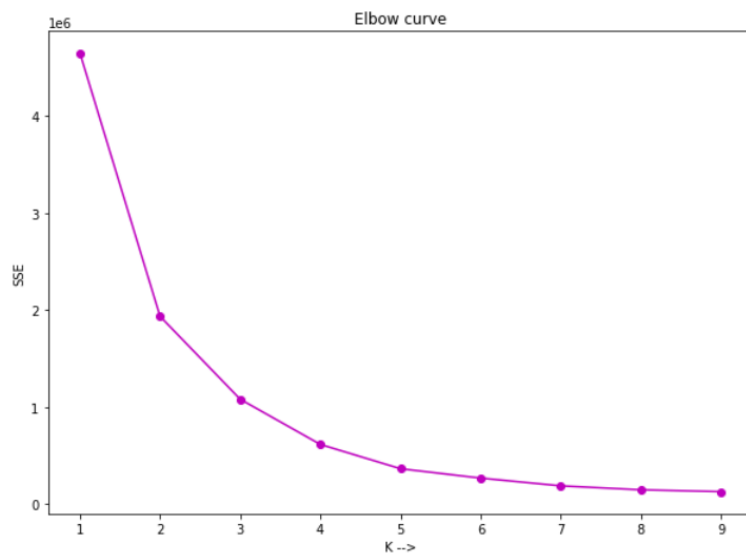
نیم سال دوم ۱۳۹۹-۱۴۰۰

سوالات پیاده سازی

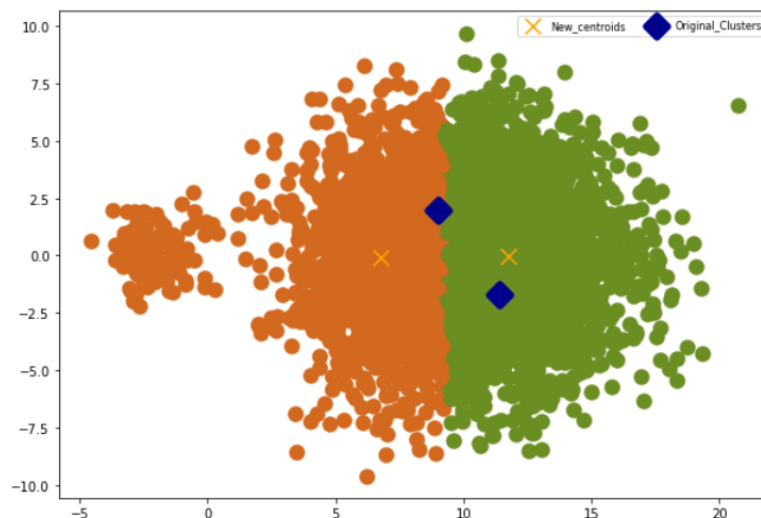
مسئله ۱:

مجموعه داده data_kmeans_1:

نمودار elbow برای مجموعه داده اول به این صورت است:

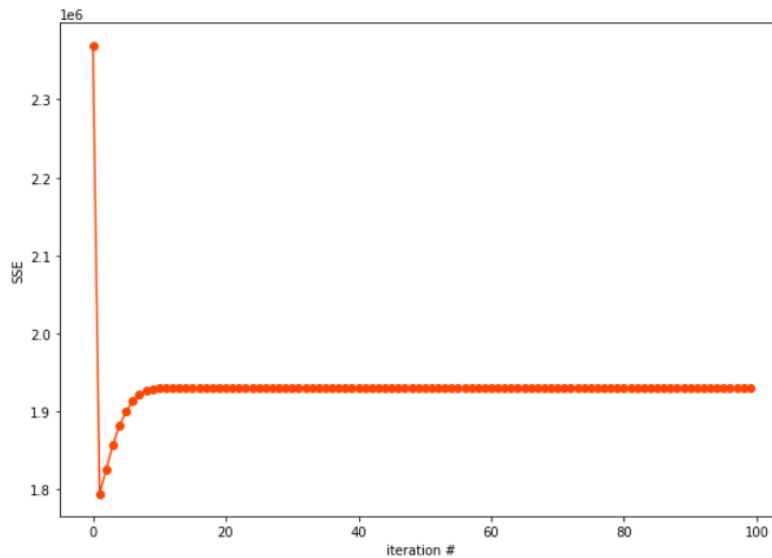


از روی نمودار به نظر می‌رسد که تعداد مناسب خوشه‌ها ۲ است. بنابراین خوشه‌بندی را با ۲ کلاستر انجام می‌دهیم:



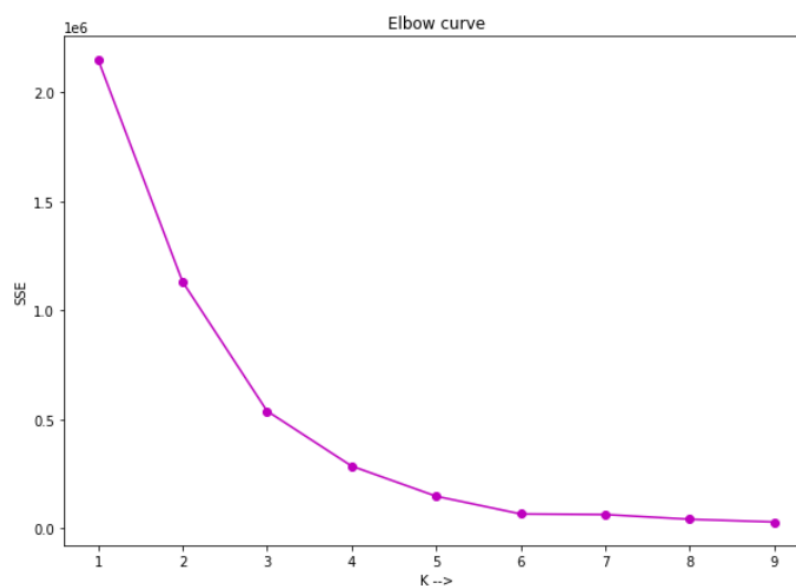
نتیجه گیری: همانطور که در نتیجه خوشه بندی برای این دیتاست مشاهده میکنید الگوریتم Kmeans وابسته به میانگین داده‌ها می‌باشد. بدین معنا که اگر یک خوشه دارای تعداد بسیاری داده نسبت به خوشه دیگر باشد، سعی میکند و بایاس است به سمت اینکه در داخل تعداد داده‌های بیشتر و یا به عبارت دیگر قسمت چگالتر داده‌ها تولید خوشه کند که مناسب نیست.

نمودار SSE به ازای تکرارهای مختلف:



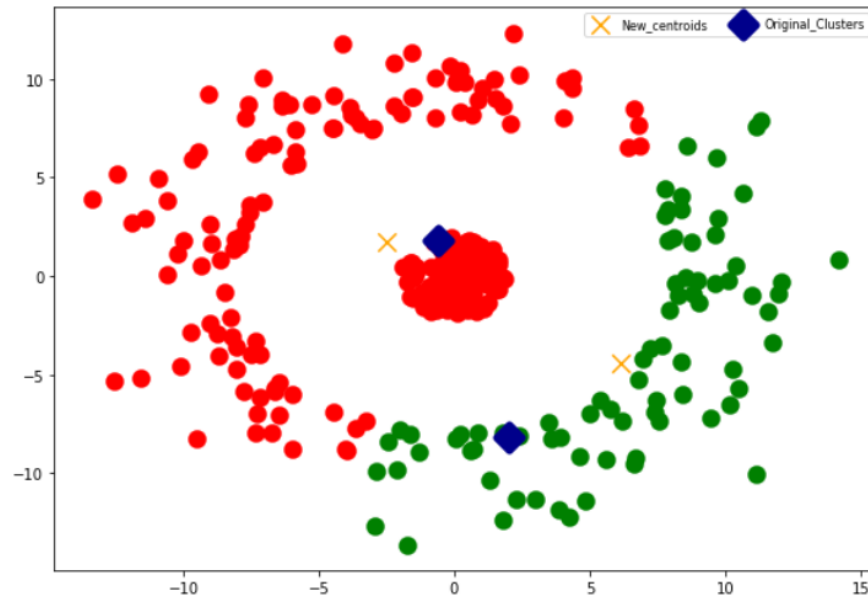
مجموعه داده data_kmeans_2 :

نمودار elbow



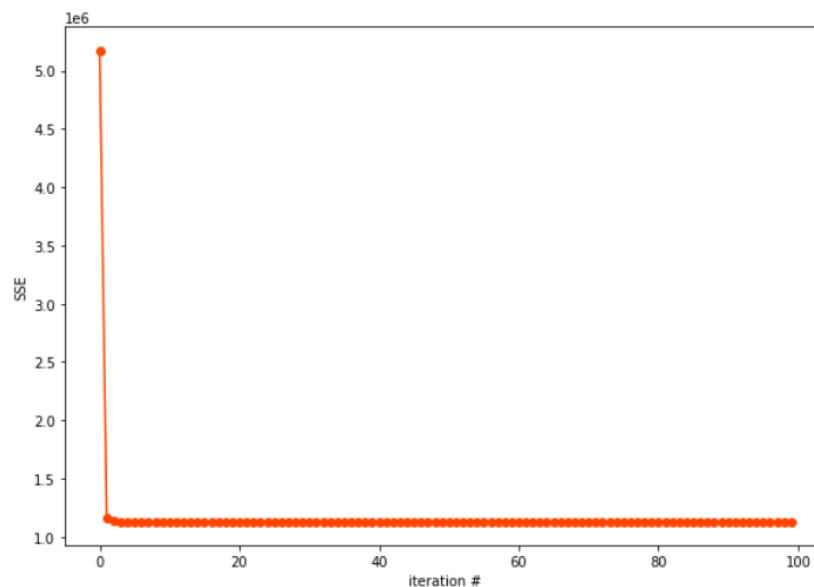
تعداد مناسب خوشه‌ها : ۲

نمودار خوشه‌بندی:

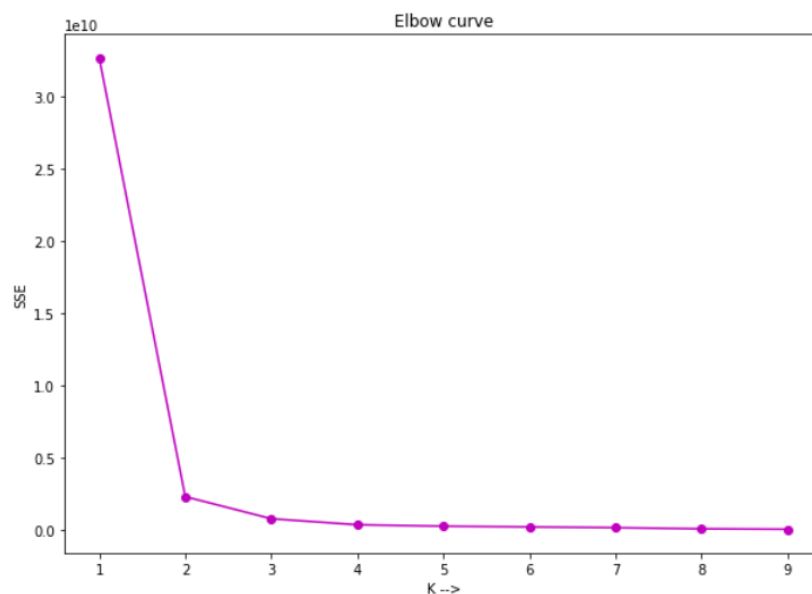


در این دادگان نیز ضعف اصلی این الگوریتم یعنی عدم توجه به چگالی و تراکم داده‌ها و صرفاً در نظر گرفتن میانگین داده‌ها به چشم می‌آید.

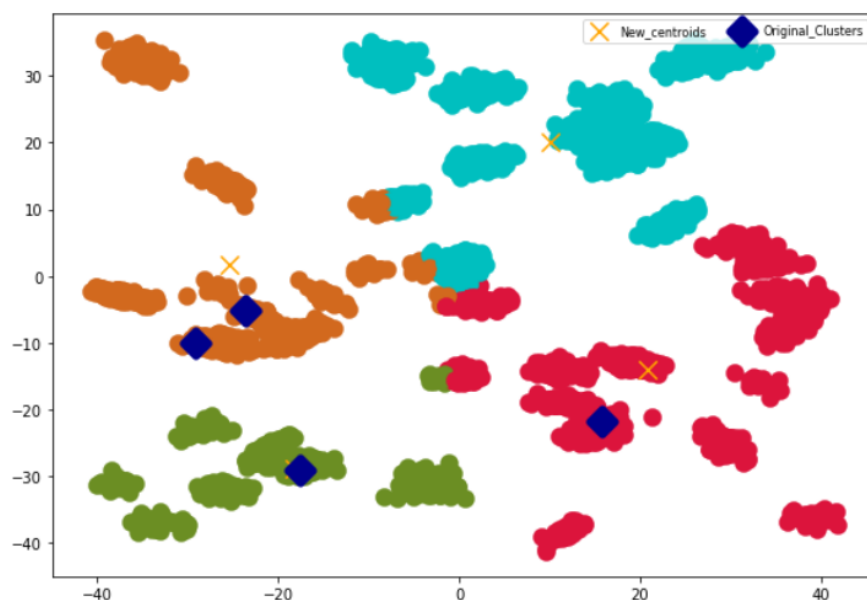
نمودار SSE به ازای تکرارهای مختلف:



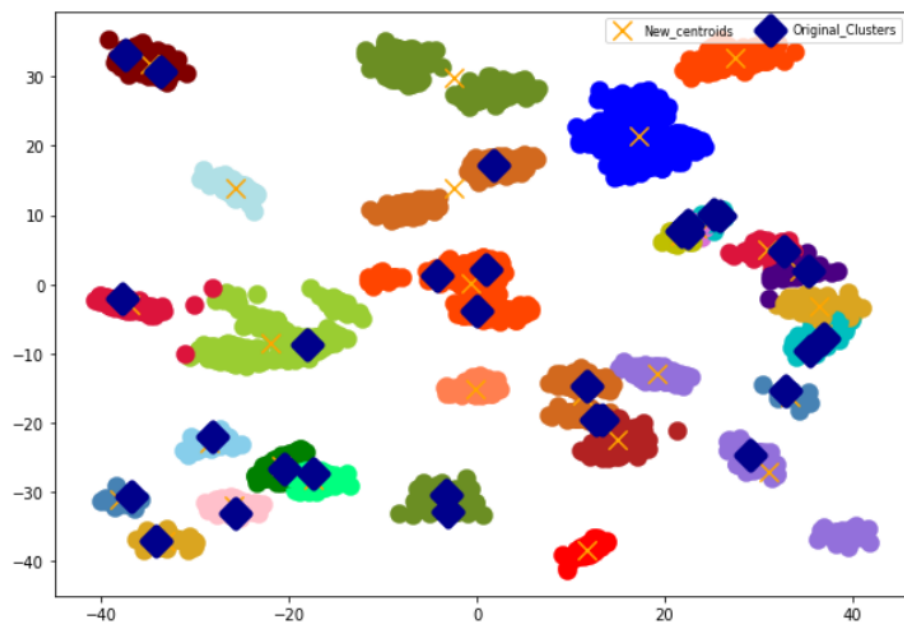
مجموعه داده data_kmeans_3 :



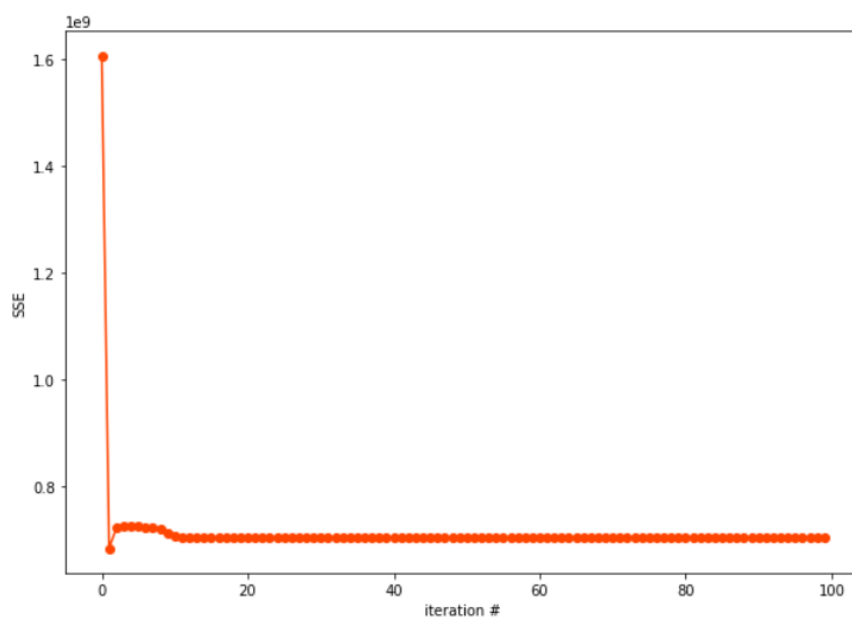
با توجه به نمودار elbow به نظر می‌رسد که تعداد مناسب خوشه‌ها حدود ۴ تا باشد. برای $k=4$ خوشه‌بندی به این صورت است:



از آنجایی که به نظر می‌رسد تعداد ۴ خوشه برای شکل داده‌ها خیلی کم است، خوشه‌بندی با $k=30$ هم انجام شد که نتیجه به صورت زیر است:

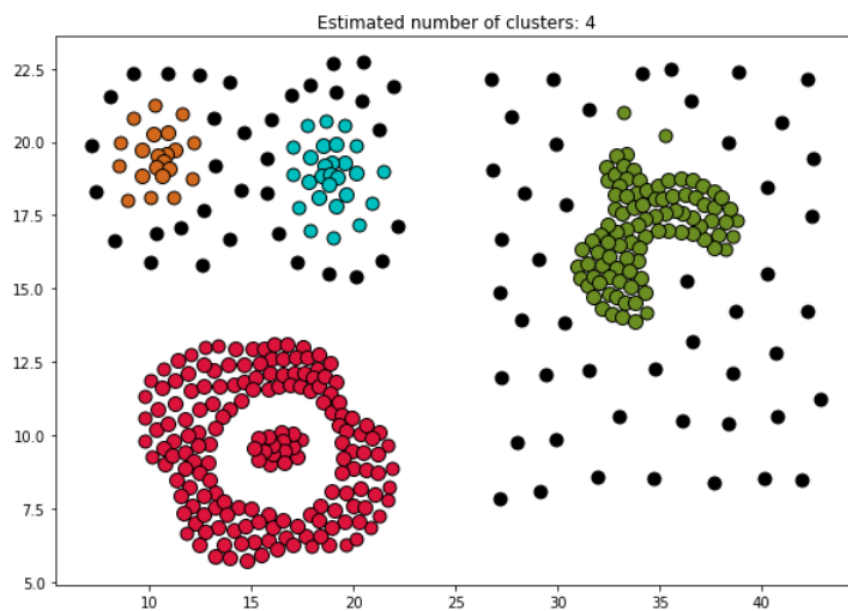
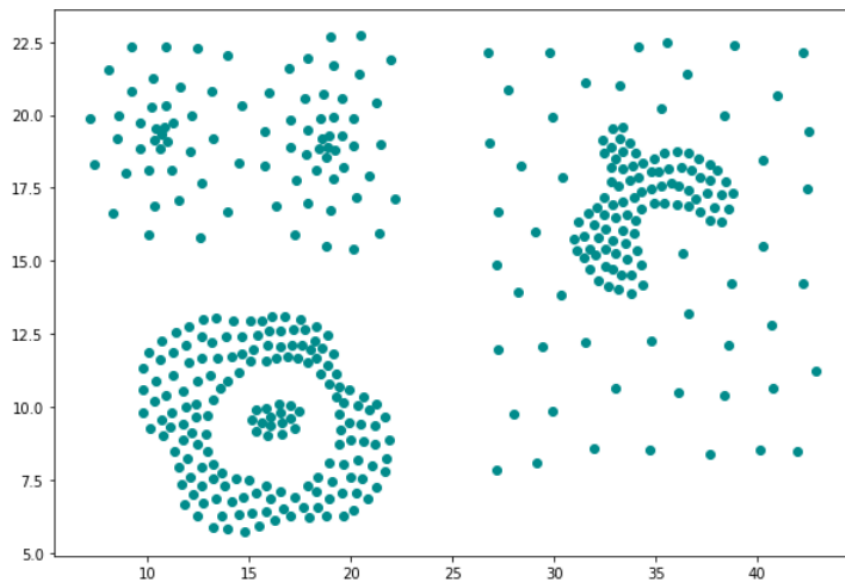


نمودار SSE به ازای تکرارهای مختلف:



مسئله ۲:

۱. مجموعه داده Compound :



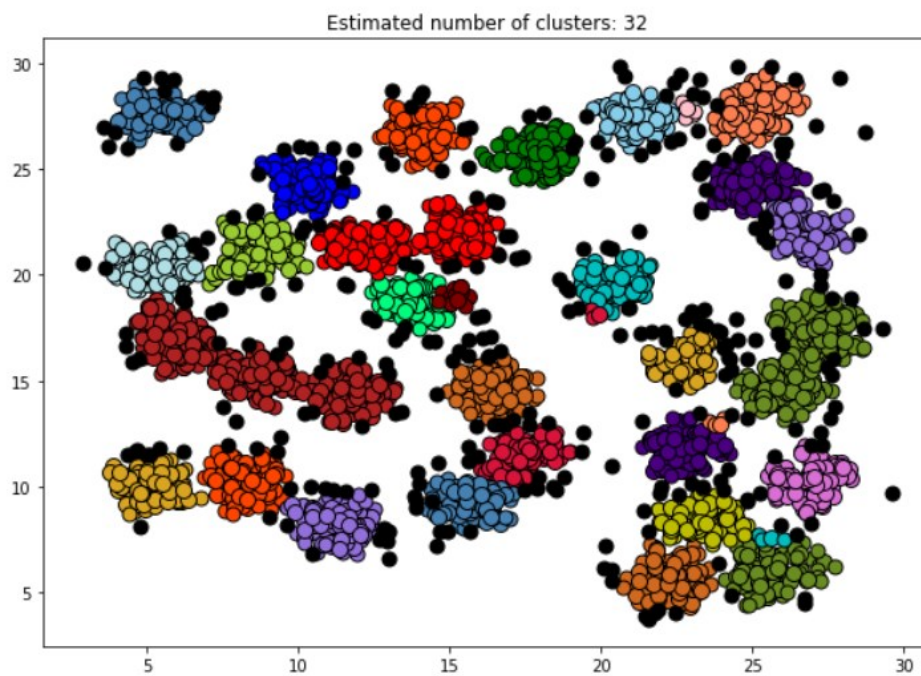
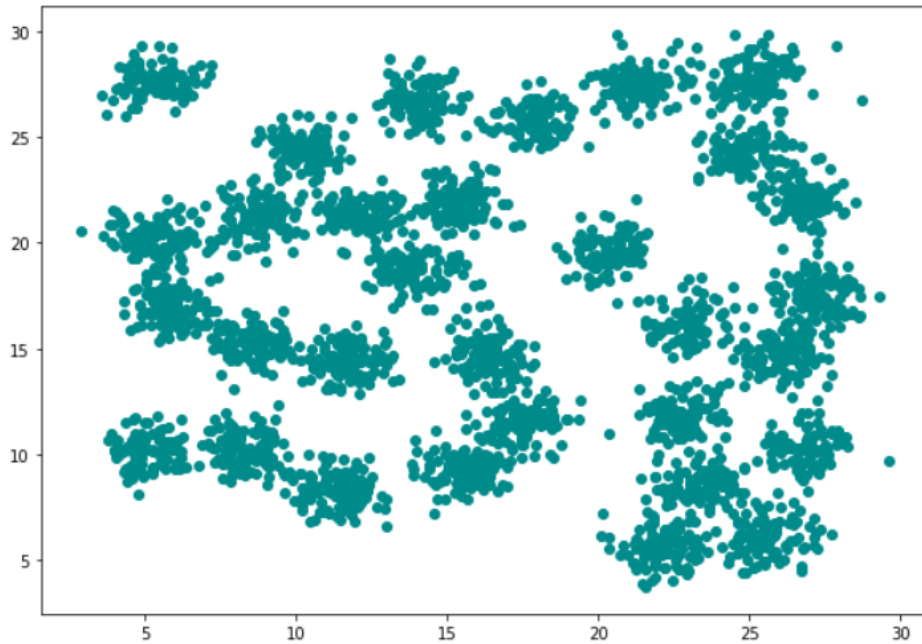
پارامترهای DBScan:

`eps=1.6, min_samples=12`

معیار purity:

Purity 0.7919799498746867

۲. مجموعه داده D31:



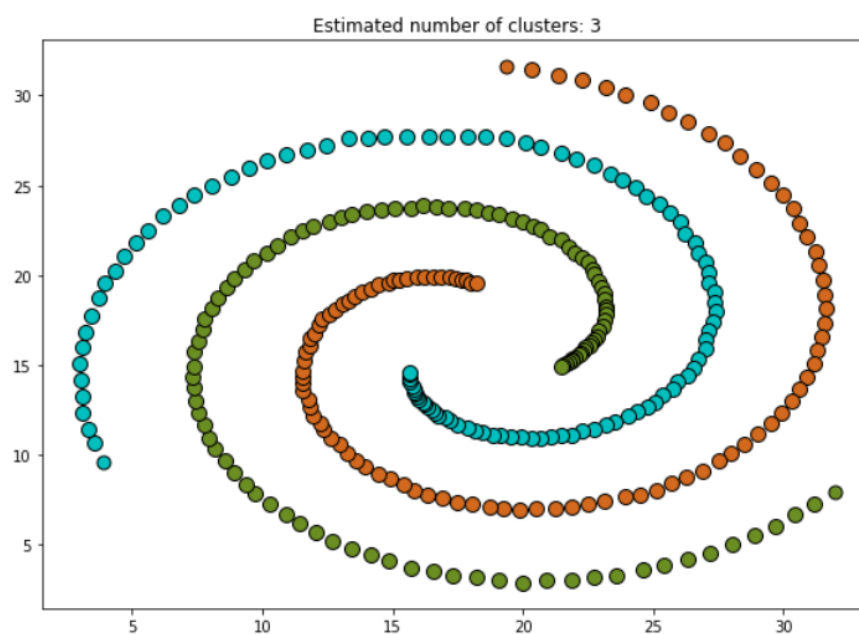
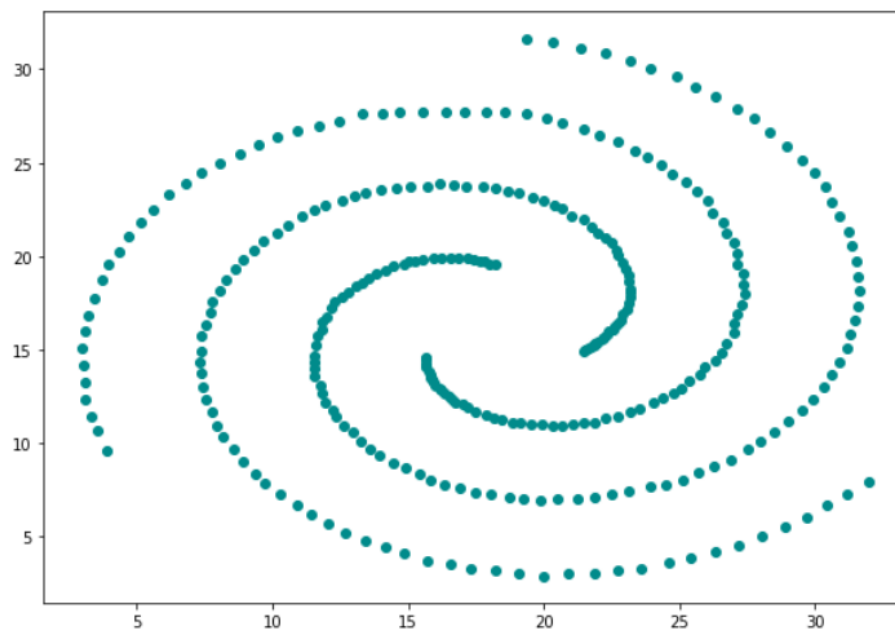
پارامترهای DBScan:

`eps=0.45, min_samples=5`

معیار purity :

Purity 0.8790322580645162

۳. مجموعه داده Spiral :



پارامترهای DBScan:

`eps=3.5, min_samples=5`

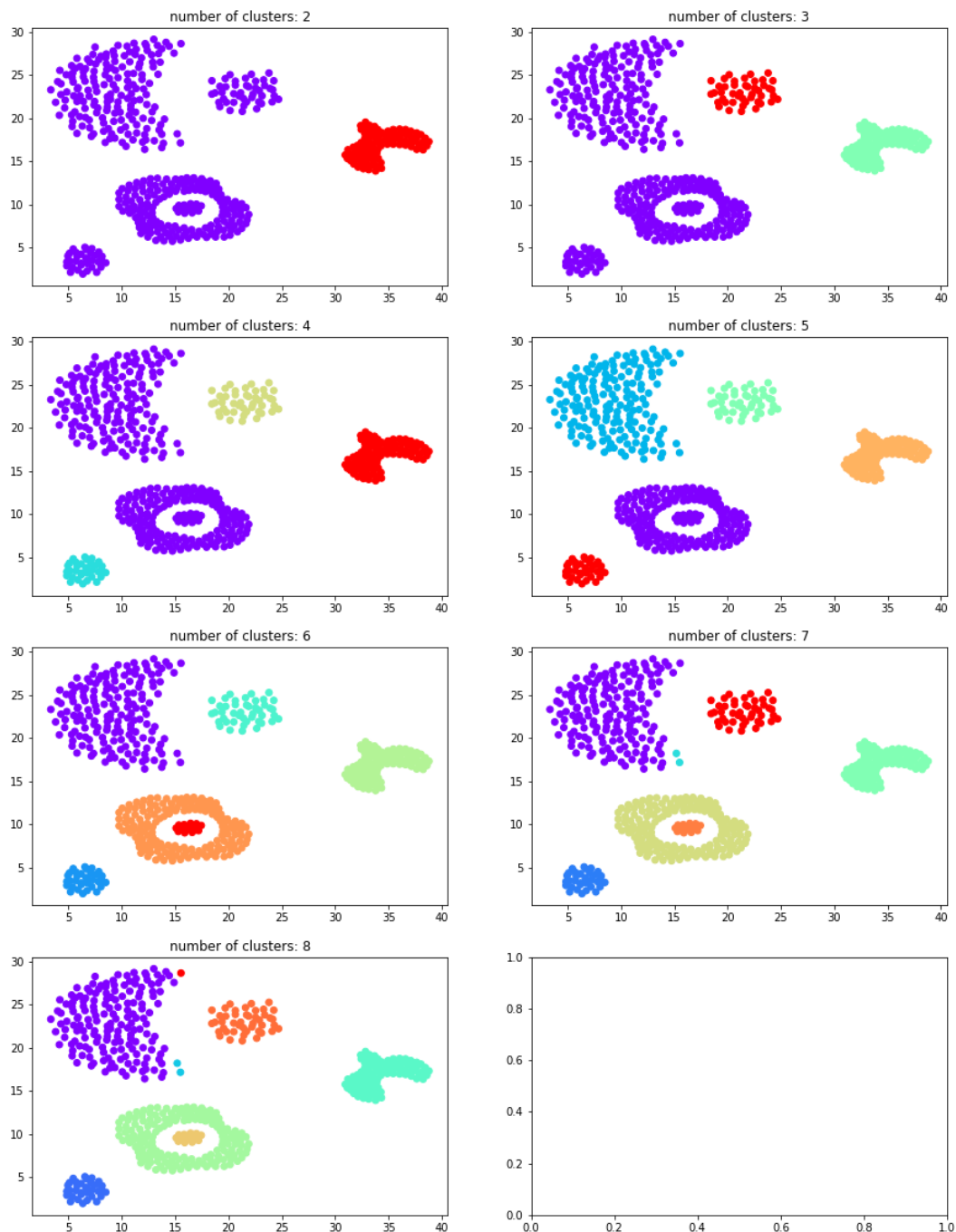
معیار purity :

Purity 1.0

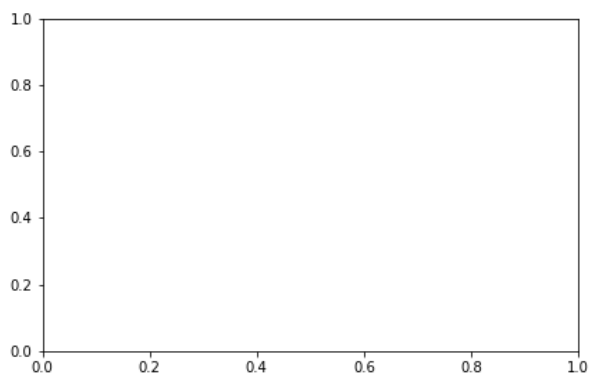
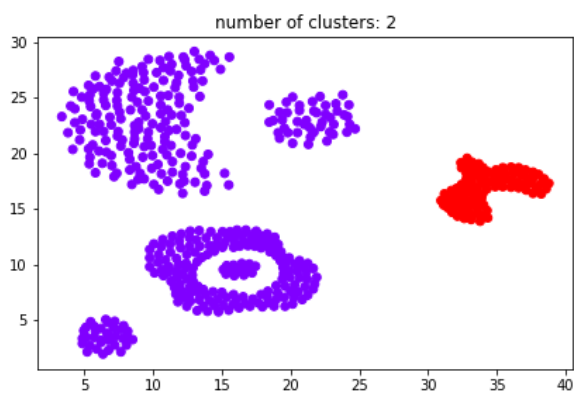
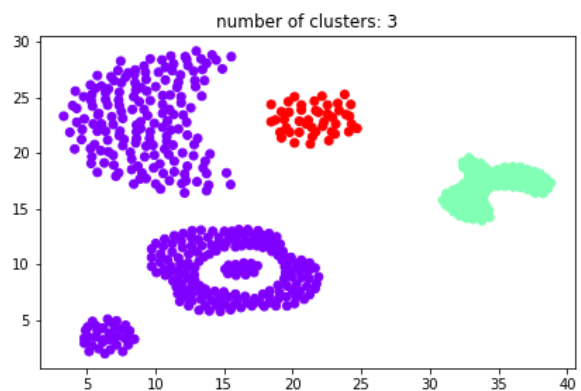
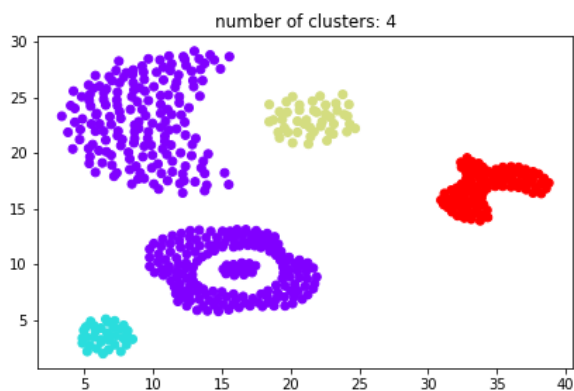
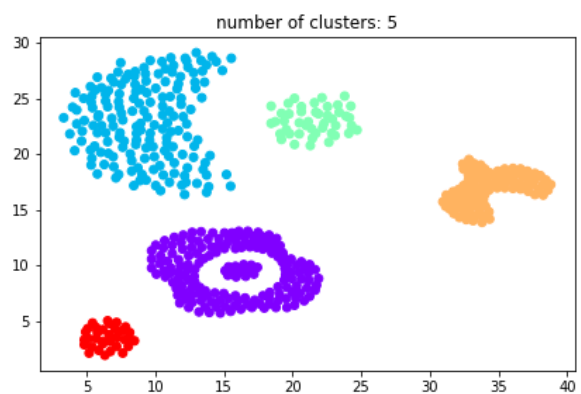
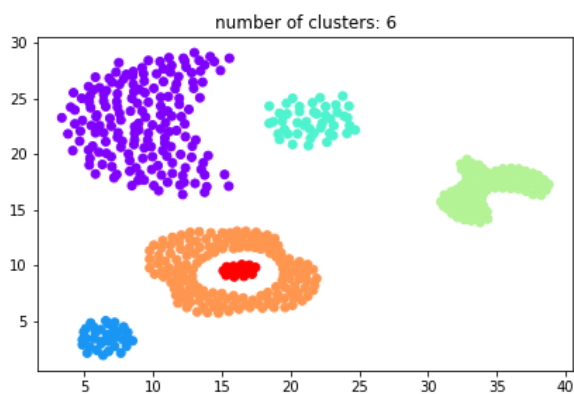
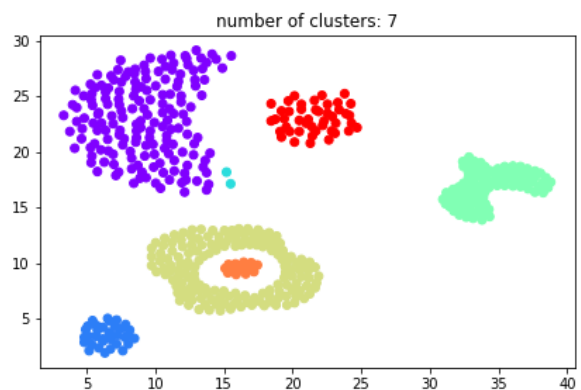
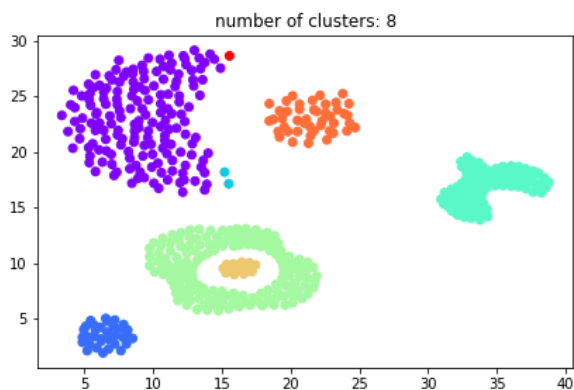
مسئله ۳:

نتایج در زیر آورده شده است :

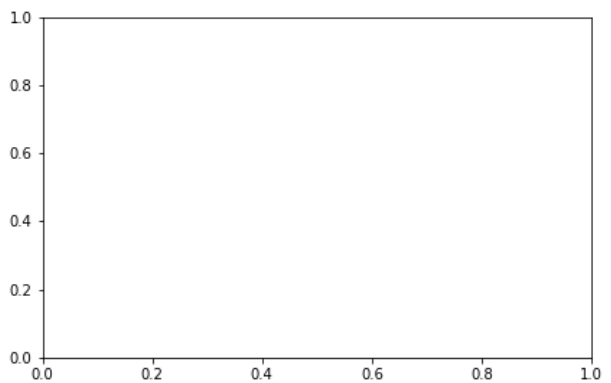
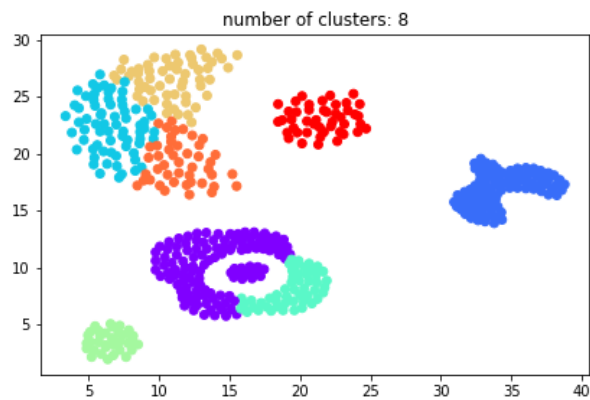
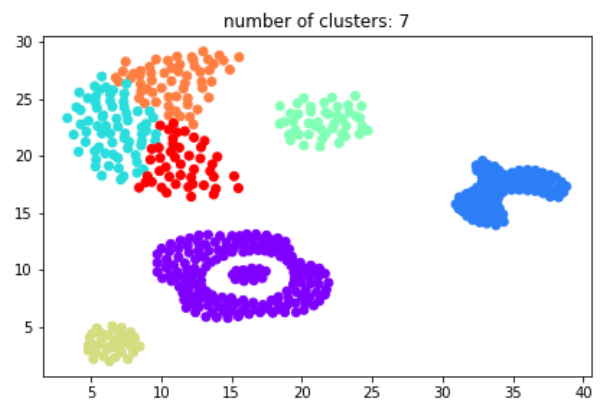
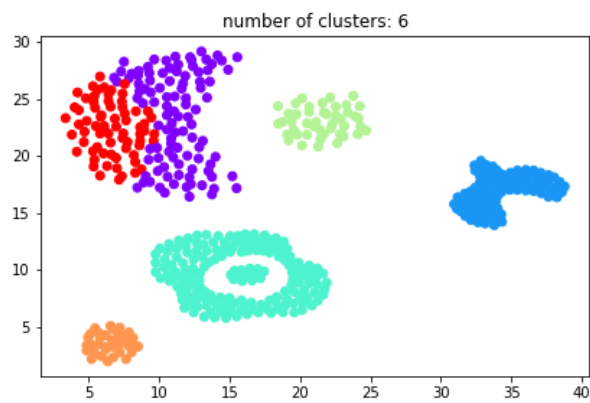
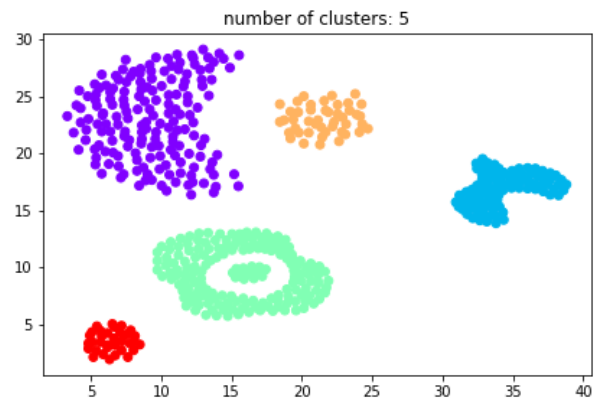
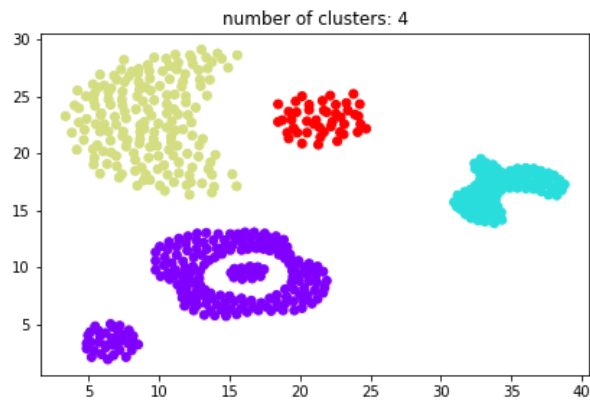
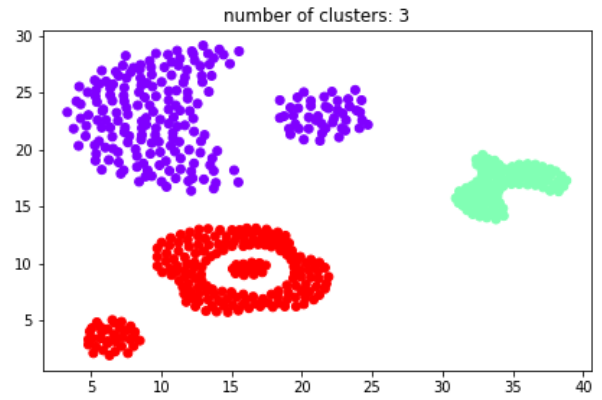
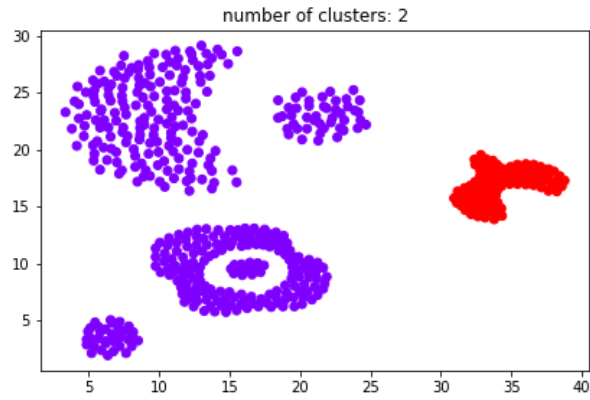
Top Down clustering - distance measure : single link



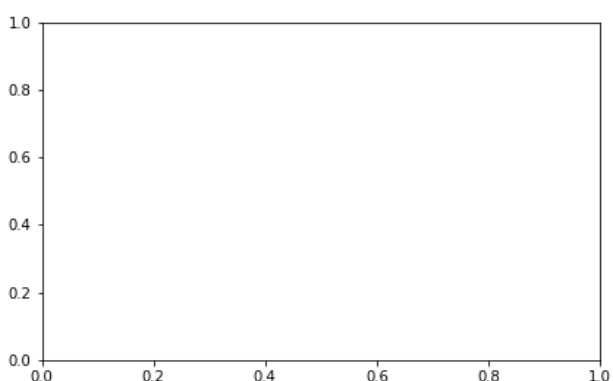
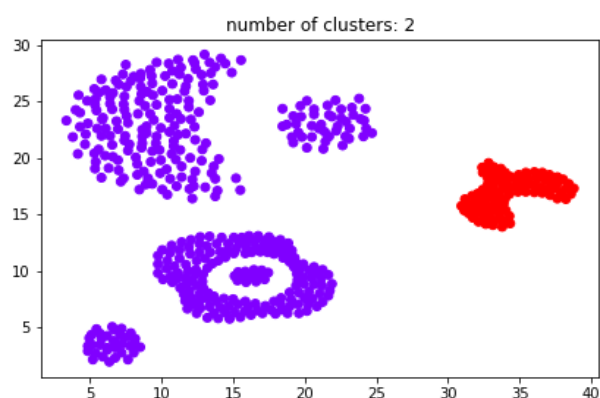
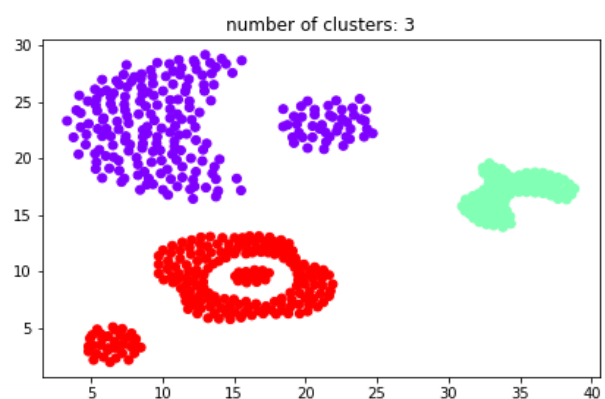
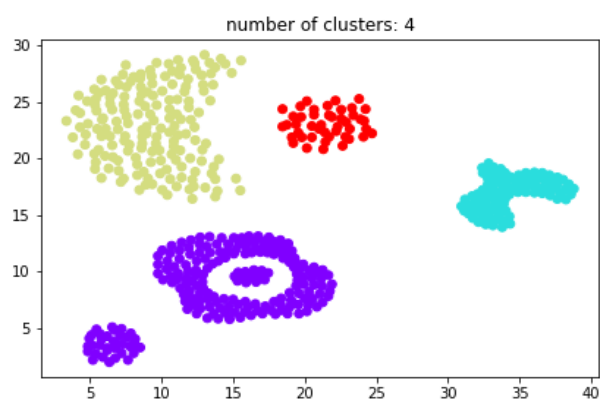
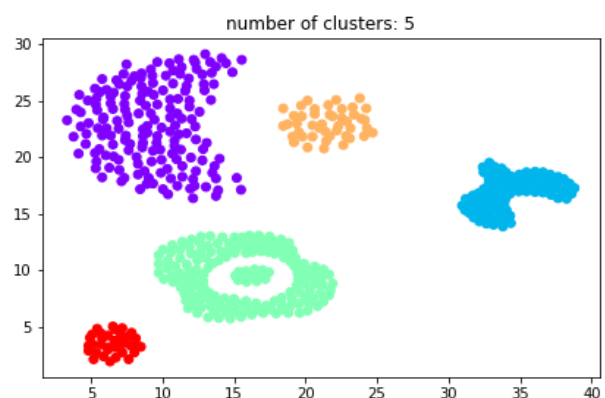
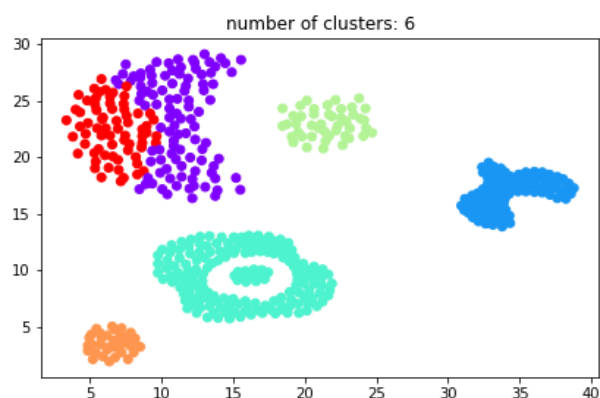
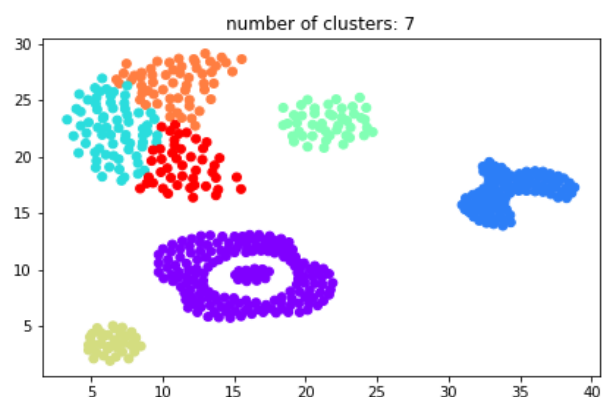
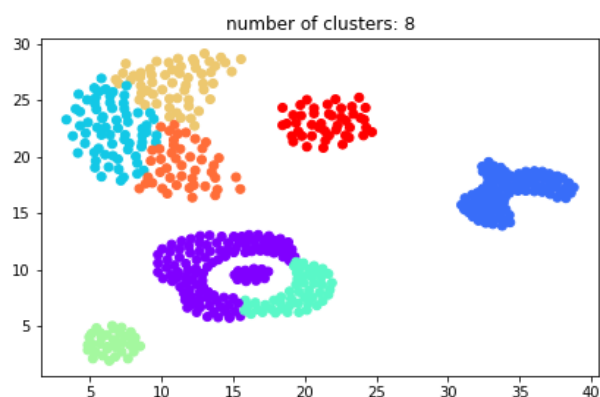
Bottom Up clustering - distance measure : single link



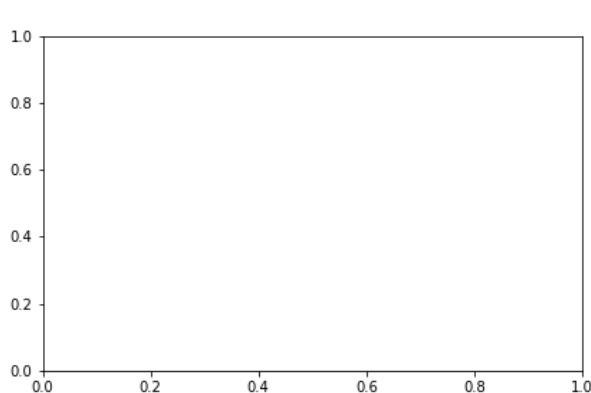
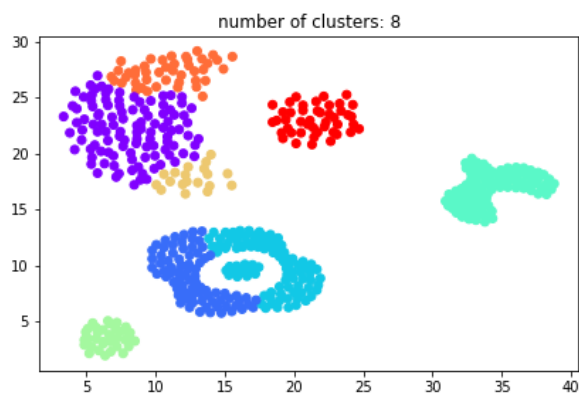
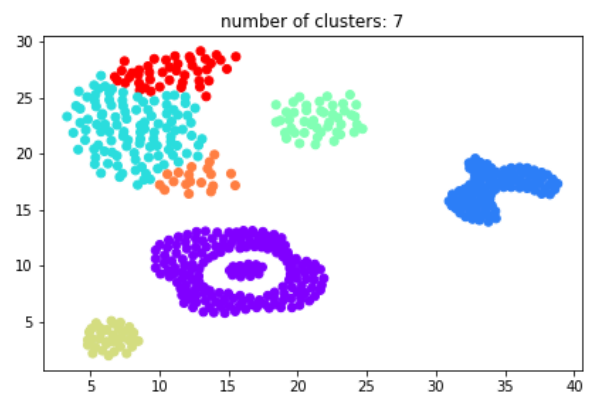
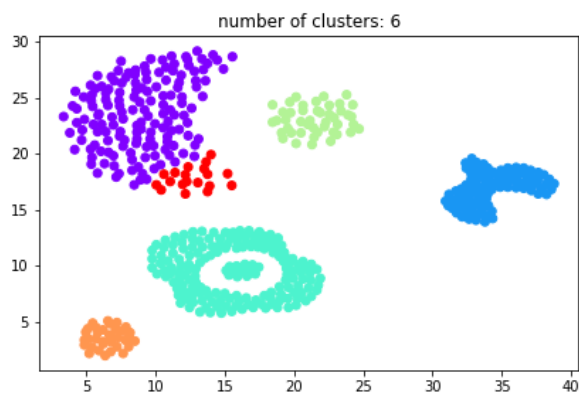
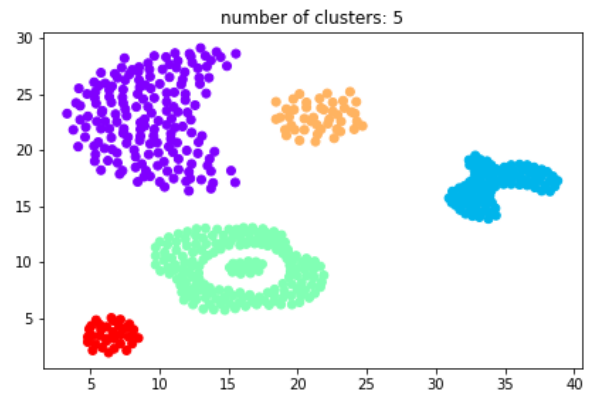
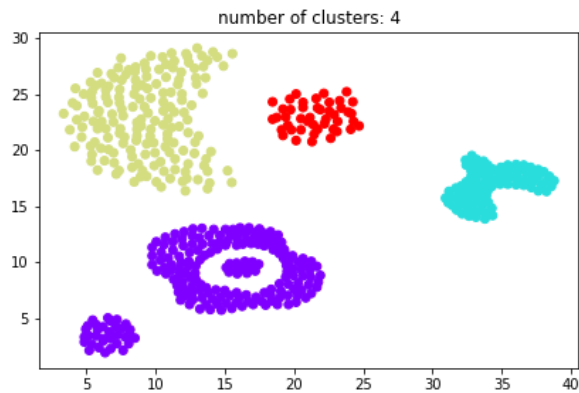
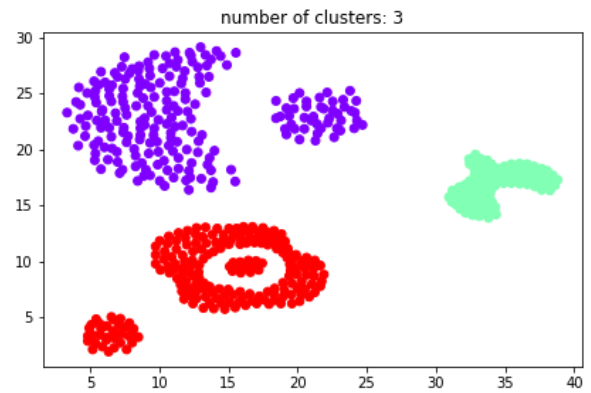
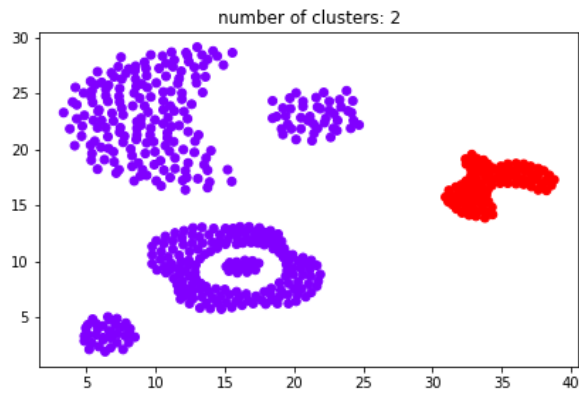
Top Down clustering - distance measure : complete link



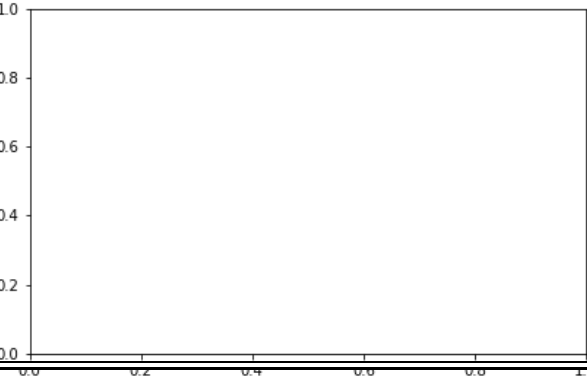
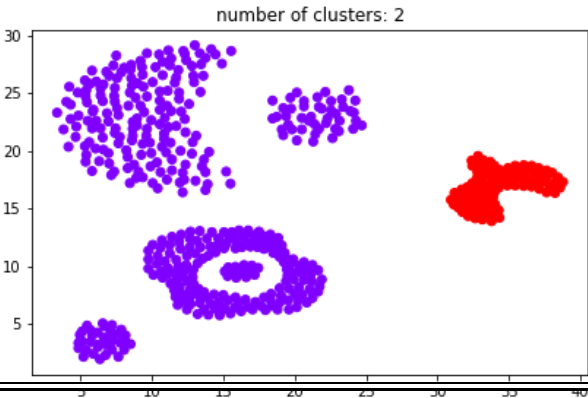
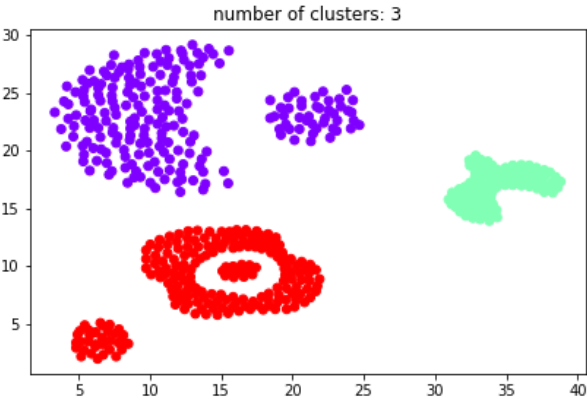
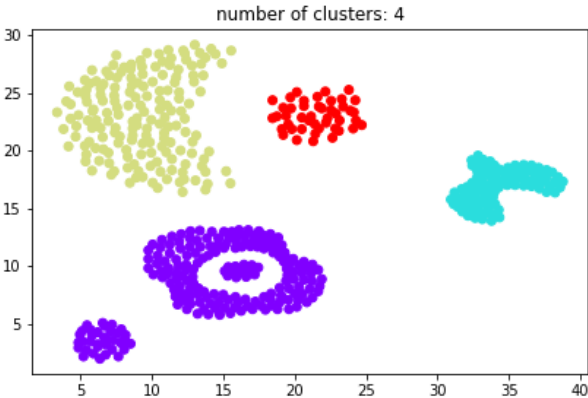
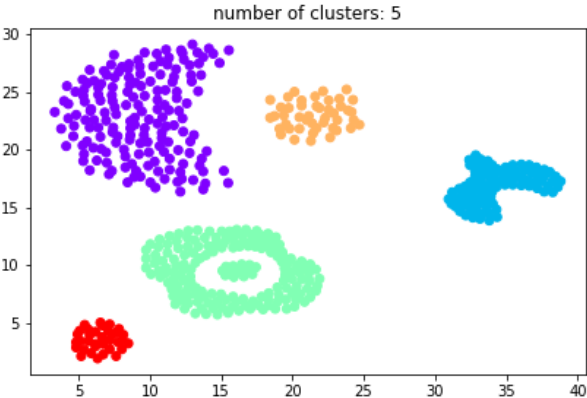
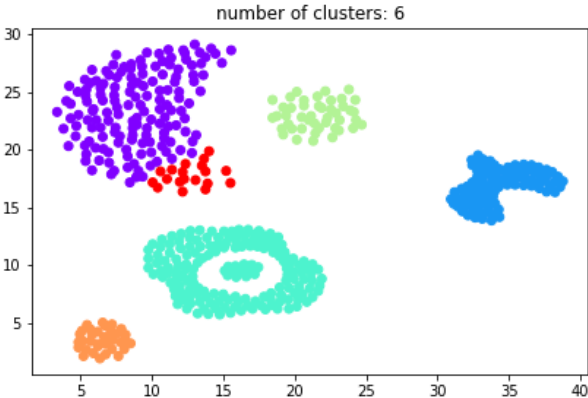
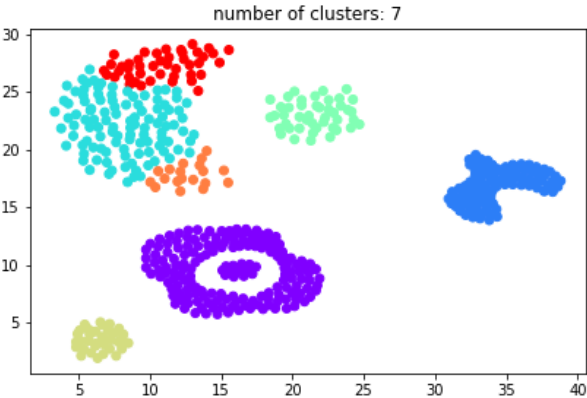
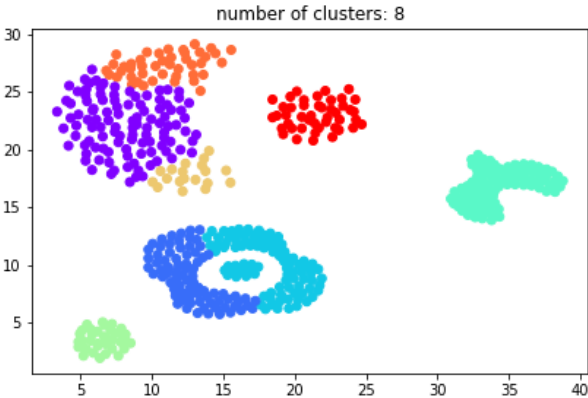
Bottom Up clustering - distance measure : complete link



Top Down clustering - distance measure : average link



Bottom Up clustering - distance measure : average link



بررسی تاثیر معیار فاصله :

در خوشه بندی سلسله مراتبی، برای محاسبه فاصله بین دو خوشه، از معیارهای AverageLink و CompleteLink.singleLink استفاده می‌شود.

دورترین فاصله یا پیوند کامل (Complete-Linkage)

$$\max\{d(a,b):a\in A,b\in B\}$$
 شیوه محاسبه

نزدیکترین فاصله یا پیوند تکی (Single-Linkage)

$$\min\{d(a,b):a\in A,b\in B\}$$
 شیوه محاسبه :

پیوند میانگین (average link)

$$(1/|A| \cdot |B|) * \sum_{a \in A} \sum_{b \in B} d(a,b)$$
 شیوه محاسبه:

از نظر حساسیت به داده پرت :

Single-link: به داده‌های پرت (outliers) حساس است. به عنوان دلیل می‌توان گفت که شاید خوشه‌هایی که به هم ماهیت نزدیکی ندارند به علت داده پرت با هم لینک شوند.

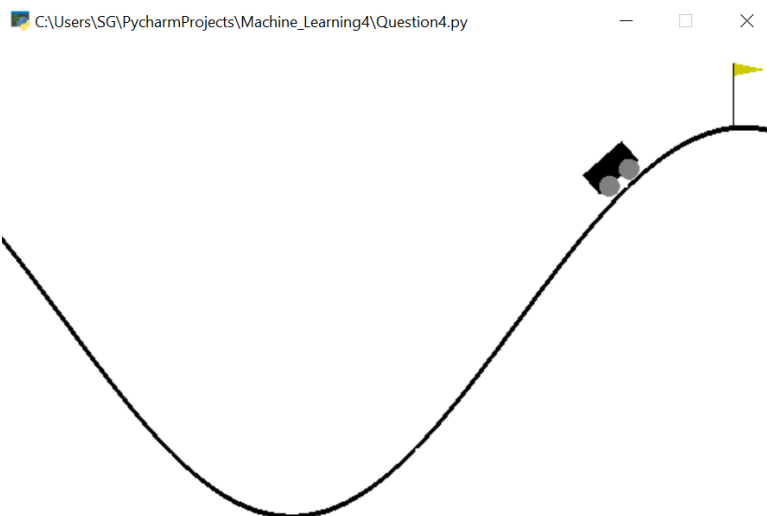
Complete-link: به شدت single-link به داده‌های پرت حساس نیست ولی باز هم حساس است. در این مورد هم می‌توان گفت شاید داده‌ها از نظر ماهیت به هم نزدیک باشند اما به علت فاصله زیاد داده‌های پرت به هم لینک نشوند.

Average-link: نسبت به دو مورد دیگر بسیار بسیار کمتر به داده‌های پرت حساسیت دارد، زیرا همواره میانگین فاصله را در نظر می‌گیرد

اگر به نتیجه خوشه‌بندی با استفاده از معیار link single نگاه کنیم مشاهده می‌کنیم که داده‌هایی تکی که به یکدیگر نزدیک هستند را به عنوان خوشه در نظر گرفته است که مناسب نیست و شاید اگر در یک دادگانی که چند داده پرت وجود داشته باشد، به مشکل شدیدتری بخوریم. از طرفی دیگر این مشکل نیز در شرایطی میتواند برای معیار link complete نیز رخ دهد و تنها به نظر میرسد معیار Average link شرایط بهتری دارد. اما در این دادگان که خیلی خاص می‌باشد به نظر بهترین نتیجه را single link و سپس link Average داشته باشد

مسئله ۴:

ابتدا باید فضای مسئله را از حالت پیوسته به حالت گسسته تبدیل کنیم. طبق کد منبع کتابخانه gym مکان اتومبیل می‌تواند در بازه [0.6 1.2] باشد و سرعت آن می‌تواند در بازه [-0.7 0.7] باشد. در ابتدا اتومبیل در نقطه ای با مختصات تصادفی ولی نزدیک به دره قرار دارد و سرعت آن صفر است. تصویری از محیط و عامل :

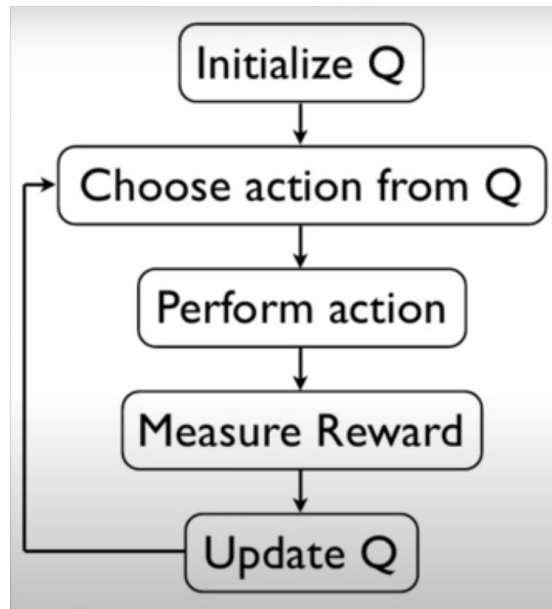


تابع Q را طبق فرمول زیر آپدیت می‌کنیم:

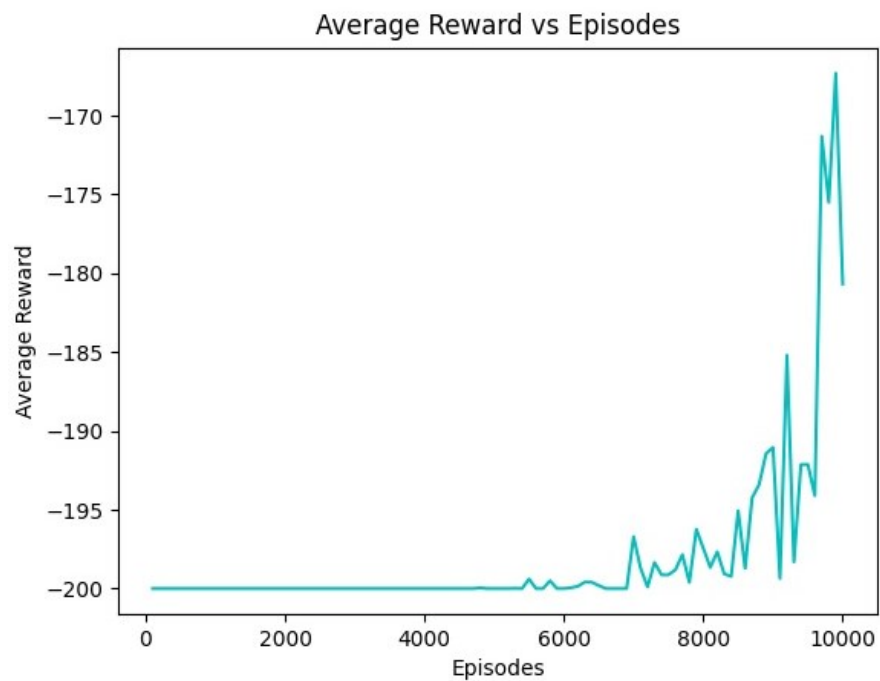
$$Q(s, a) = R(s, a) + \gamma \sum_{s'} \left(P(s, a, s') \max_{a'} Q(s', a') \right)$$

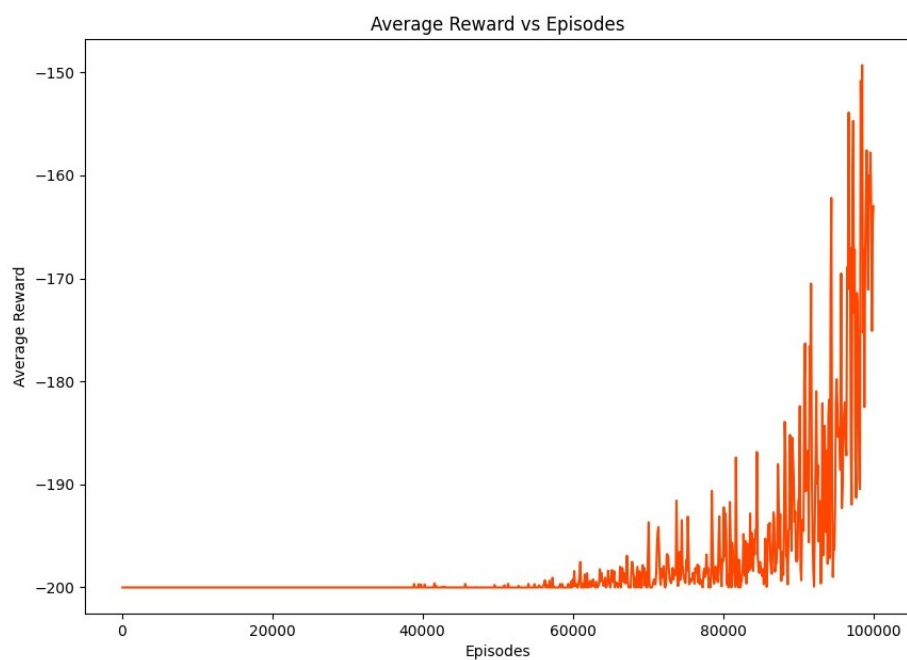
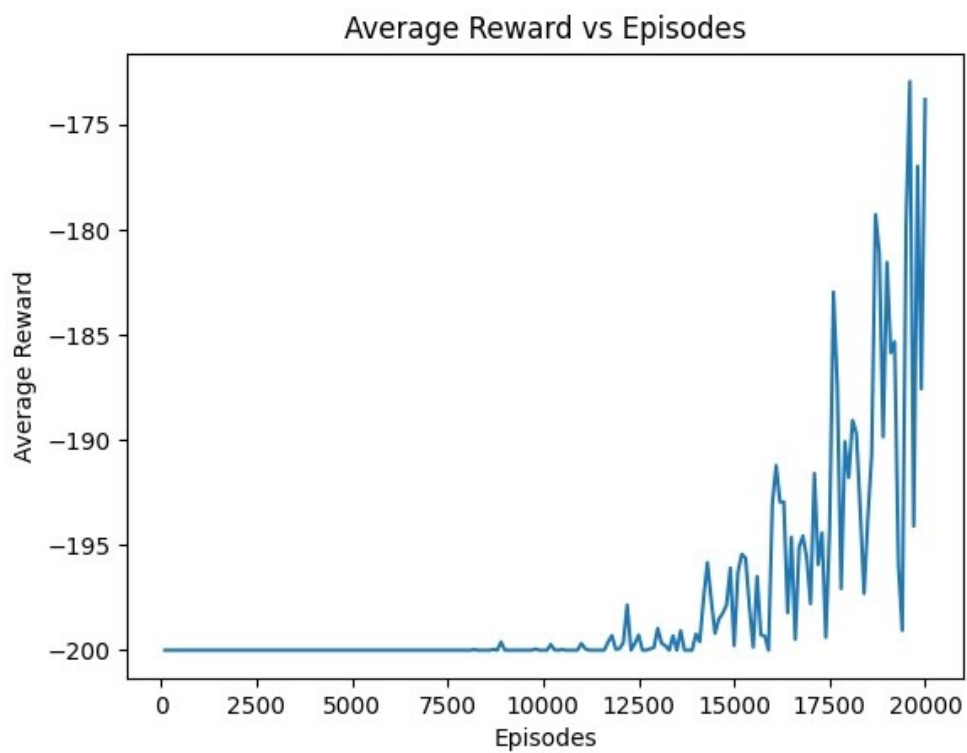
برای موازنه میان exploration و exploitation به طور تصادفی و با احتمال مشخصی برای قدم بعدی، بین حالت بهینه و حالت رندوم یکی را انتخاب می‌کنیم. همچنین احتمال انتخاب حالت حریصانه را با گذشت زمان زیاد می‌کنیم.

فلوچارت کلی الگوریتم در شکل زیر آورده شده است:



نمودار میانگین پاداش دریافتی عامل به ازای تعداد مختلفی از تکرارها در شکل‌های زیر رسم شده است:





همانطور که مشاهده می‌شود با اینکه نمودارها حالت نویزی دارند اما در همه آن‌ها استراتژی کلی افزایش پاداش است یعنی عامل سعی دارد که در دراز مدت پاداش خود را ماکزیمم کند.