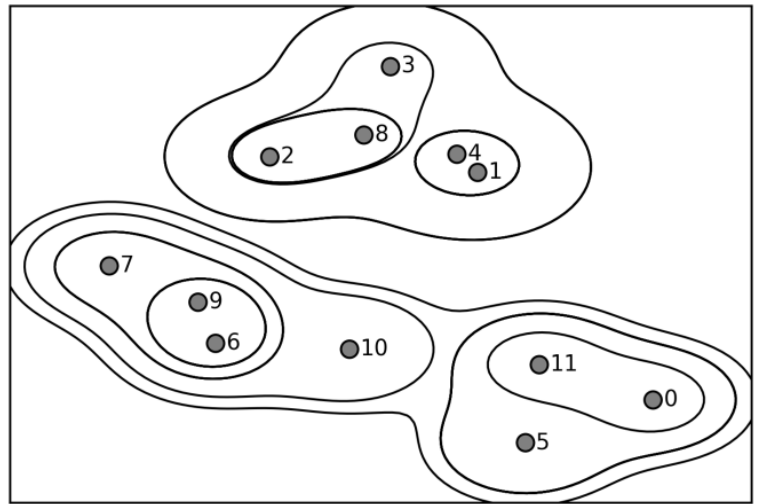


<p>• نقطه پرت^۶ نقطه‌ای است که نقطه هسته نبوده و متعلق به هیچ خوشه‌ای نباشند. از آنجایی که در الگوریتم DBScan هر نقطه هسته، الزاماً تشکیل‌دهنده یک خوشه است می‌توان تعریف نقطه پرت را به این صورت بیان کرد که نقطه پرت هسته نیست و در دسترس چگال از یک نقطه هسته دیگر قرار ندارد. اگر نقطه‌ای، نقطه هسته نبوده اما در دسترس چگال از یک نقطه هسته باشد، نقطه مرزی^۷ نامیده می‌شود.</p> <p>حال با مفاهیم فوق می‌توان الگوریتم DBScan را بیان کرد. این الگوریتم هر نقطه را با توجه به محدوده آن بررسی می‌کند. هرکدام از نقطه‌ها در یکی از سه گروه هسته، مرزی یا پرت قرار می‌گیرند. در ادامه الگوریتم، هسته‌هایی که در دسترس چگال هم هستند، ادغام‌شده و خوشه واحدی می‌شوند. با این اوصاف هر خوشه از حداقل یک هسته و تعدادی نقاط مرزی تشکیل می‌شود.</p> <pre>from sklearn.cluster import DBSCAN</pre> <pre>...</pre> <pre>X = iris_data[['sepal.length', 'sepal.width', 'petal.length', 'petal.width']]</pre> <pre>dbscan = DBSCAN(eps=0.5, min_samples=5)</pre> <pre>labels = dbscan.fit_predict(X)</pre> <pre>noise = X[labels == -1]</pre> <p>خوشه‌بندی سلسله‌مراتبی یکی از روش‌های گروه‌بندی داده‌ها است که بر اساس ساختار دسته‌بندی تدریجی کار می‌کند. این روش داده‌ها را به صورت ساختار درختی مرتب می‌کند تا ارتباط بین گروه‌های مختلف را مشخص کند.</p> <p>شروع با داده‌های مجزا: ۱. ابتدا هر داده به عنوان یک خوشه‌ی مستقل در نظر گرفته می‌شود. ۲. ادغام خوشه‌های مشابه: در هر مرحله، دو خوشه‌ای که بیشترین شباهت را دارند، با یکدیگر ترکیب می‌شوند. ۳. ساختار سلسله‌مراتبی: این روند ادامه می‌یابد تا ...</p> <pre>from sklearn.cluster import AgglomerativeClustering</pre> <pre>clustering = AgglomerativeClustering(n_clusters=3)</pre> <pre>labels = clustering.fit_predict(X)</pre>	<p>الگوریتم خوشه‌بندی DBScan، از جمله الگوریتم‌های خوشه‌بندی معروف است. نام این الگوریتم که برگرفته از حروف اول عبارت <u>Density Based Spatial Clustering of Applications with Noise</u> می‌باشد، بر اساس چگالی نقاط موجود در اطراف هر نقطه پایه‌گذاری شده و خصوصاً برای یافتن خوشه‌ها در مجموعه داده‌هایی با خوشه‌هایی به اشکال دلخواه ارائه شده است. این الگوریتم نیازمند دو پارامتر Eps و MinPts است. قابل ذکر است که Eps یک عدد حقیقی مثبت که بیان‌کننده اندازه شعاع برای محاسبه چگالی نقاط موجود در اطراف هر نقطه به اندازه آن پارامتر می‌باشد. پارامتر MinPts نیز یک عدد طبیعی است که اصولاً بزرگ‌تر از ۱ انتخاب می‌شود. الگوریتم DBScan با بهره‌گیری از تعاریف و مفاهیم زیر که بر مبنای پارامترهای Eps و MinPts هستند، طراحی می‌شود.</p> <ul style="list-style-type: none"> • همسایگی Eps یک نقطه^۱: همسایگی Eps یک نقطه مانند p شامل مجموعه نقاطی مانند q از مجموعه داده‌ای است طوری که فاصله q از p بیشتر از Eps نباشد. • نقاط هسته^۲: نقطه q، یک نقطه هسته نامیده می‌شود اگر تعداد نقاط موجود در محدوده به شعاع Eps آن با احتساب خود q کمتر از MinPts نباشد. • مجموعه نقاط محدود: یک نقطه: نقطه p در مجموعه محدود نقاط q یا به اختصار محدود q است اگر فاصله آن تا q بیشتر از Eps نباشد. • قابل دسترس چگال به طرز مستقیم یا اختصاراً دسترس چگال مستقیم^۳: نقطه p در دسترس چگال مستقیم q است اگر نقطه q هسته بوده و p نیز در محدوده q باشد. • دسترس چگال^۴: نقطه p در دسترس چگال q است اگر زنجیره p_1, p_2, \dots, p_n موجود باشد که $p_1 = q$، $p_n = p$ و p_{i+1} در دسترس چگال مستقیم p_i باشد. • پیوستگی چگال^۵: نقطه p با q پیوستگی چگال دارد اگر نقطه‌ای مانند o موجود باشد، طوری که p و q در دسترس چگال o باشند. • خوشه: زیرمجموعه C از مجموعه داده‌ای D یک خوشه محسوب می‌شود اگر: <ol style="list-style-type: none"> ۱. به از هر دونقطه مانند p و q از مجموعه داده‌ای D، اگر p عضو C بوده و q در دسترس چگال p باشد در این صورت q نیز عضو C باشد. ۲. به از هر دونقطه مانند p و q از زیرمجموعه C، p باید با q پیوستگی چگال داشته باشد.
---	--

¹ Eps-neighborhood of a point² Core point³ Directly density-reachable⁴ Density-reachable⁵ Density-connected⁶ Noise⁷ Boarder point

تی اس ان ای (t-sne) و پی سی ای (PCA) هر دو از تکنیک‌های کاهش ابعاد هستند، اما با اهداف و روش‌های متفاوت:

- پی سی ای (PCA) یک روش خطی است که برای کاهش ابعاد داده‌های دارای ساختار خطی استفاده می‌شود. این روش داده‌ها را به محورهای جدید تبدیل می‌کند که بیشترین واریانس را دارند، یعنی اطلاعات اصلی را حفظ می‌کند و نویز را کاهش می‌دهد.
 - تی اس ان ای (t-SNE) یک روش غیرخطی است که عمدتاً برای تصویری‌سازی داده‌های پیچیده استفاده می‌شود. این روش سعی می‌کند نقاط داده‌ای مشابه را به یکدیگر نزدیک کند، در حالی که فاصله‌های بین داده‌های غیرمشابه را حفظ می‌کند، که باعث ایجاد ساختار خوشه‌ای واضح‌تر در داده‌های تصویری شده می‌شود.
- به طور خلاصه، پی سی ای برای کاهش ابعاد کارآمد و سریع است، ولی تی اس ان ای بیشتر برای نمایش و درک خوشه‌ها در داده‌های پیچیده کاربرد دارد.



تمرین - برنامه زیر را اجرا کنید و نتایج را گزارش دهید

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.cluster.hierarchy import dendrogram, linkage
from sklearn.datasets import make_blobs
from sklearn.cluster import AgglomerativeClustering
```

لود آیریس

X = ...

Apply hierarchical clustering

linkage_matrix = linkage(X, method='ward')

Plot dendrogram

plt.figure(figsize=(10, 5))

dendrogram(linkage_matrix)

plt.title("Dendrogram - Hierarchical Clustering")

plt.xlabel("Sample Index")

plt.ylabel("Distance")

plt.show()

Perform clustering

clustering = AgglomerativeClustering(n_clusters=3)

labels = clustering.fit_predict(X)

Plot clusters

plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis', s=50, edgecolors='k')

plt.title("Hierarchical Clustering Results")

plt.xlabel("Feature 1")

plt.ylabel("Feature 2")

plt.show()