

سوال ۳: یکی از ایده های اصلی در کلاسترینگ استفاده از فاصله بین نقاط است. آیا این روش همیشه جواب می دهد؟ در چه شرایطی این روش می تواند نتیجه منفی بدهد.

روش های مبتنی بر فاصله در کلاسترینگ (مانند K-Means یا Hierarchical Clustering) از فاصله بین نقاط برای تعیین شباهت یا نزدیکی استفاده می کنند. اگرچه این روش ها در بسیاری از سناریوها مفید هستند، اما در شرایط خاصی ممکن است به نتایج نامطلوب منجر شوند. برخی از این شرایط:

۱. داده های غیرمحدب یا پیچیده: روش های مبتنی بر فاصله برای کلاسترهایی که به صورت محدب هستند (مانند دایره یا بیضی) عملکرد خوبی دارند. اما اگر داده ها شکل های پیچیده ای داشته باشند (مثل کلاسترهای دایره ای داخل یکدیگر یا داده های خطی پیچیده)، این روش ها ممکن است نتوانند کلاسترها را به درستی تشخیص دهند.

۲. مقیاس گذاری نادرست ویژگی ها: اگر ویژگی های داده ها در مقیاس های مختلف باشند و بدون نرمال سازی یا استانداردسازی استفاده شوند، فاصله ها به طور نادرست محاسبه می شوند. این می تواند باعث شود که ویژگی هایی با مقیاس بزرگتر، اهمیت بیشتری پیدا کنند.

۳. داده های نویزی یا با نویز بالا: داده هایی که دارای نویز یا نقاط پرت (outliers) زیادی هستند می توانند محاسبات فاصله را مختل کنند، زیرا این نقاط فاصله های غیرمعمولی ایجاد می کنند که تأثیر منفی روی الگوریتم دارند.

۴. ابعاد بالا (مشکل Curse of Dimensionality): در داده های با ابعاد بالا، فاصله بین نقاط تمایل دارد یکنواخت شود (یعنی تفاوت فاصله بین نزدیک ترین و دورترین نقاط کاهش می یابد). این باعث می شود که الگوریتم های مبتنی بر فاصله قدرت خود را از دست بدهند و نتوانند کلاسترها را به خوبی تشخیص دهند.

۵. فرضیات نادرست درباره تعداد کلاسترها: بسیاری از الگوریتم های مبتنی بر فاصله (مانند K-Means) نیاز به تعداد کلاسترها به عنوان ورودی دارند. اگر تعداد کلاسترها به درستی مشخص نشود، نتایج کلاسترینگ ممکن است بی معنی باشد.

۶. همگن نبودن چگالی کلاسترها: اگر کلاسترها چگالی های متفاوتی داشته باشند (مثلاً برخی کلاسترها دارای نقاط متراکم و برخی دارای نقاط پراکنده باشند)، روش های مبتنی بر فاصله ممکن است نتوانند این کلاسترها را به درستی شناسایی کنند.

الگوریتم DBSCAN را توضیح دهید و همچنین توضیح دهید در کدام دسته از الگوریتم های کلاسترینگ قرار می گیرد. در ادامه تفاوت آن را با الگوریتم OPTICS توضیح دهید.

DBSCAN یکی از الگوریتم های کلاسترینگ مبتنی بر چگالی است که برای پیدا کردن کلاسترها در داده هایی که ساختار نامنظم دارند یا شامل نویز هستند، مناسب است. پارامترها:

- $\epsilon$ : شعاعی که همسایگی یک نقطه را مشخص می کند.
- $minPts$ : حداقل تعداد نقاط مورد نیاز برای تشکیل یک کلاستر در یک منطقه  $\epsilon$ -همسایگی.

تعریف نقاط:

- نقطه مرکزی (Core Point): نقطه‌ای که حداقل  $minPts$  نقطه (شامل خودش) در شعاع  $\epsilon$  آن وجود داشته باشد.

- نقطه مرزی (Border Point): نقطه‌ای که در شعاع  $\epsilon$  یک نقطه مرکزی قرار دارد اما خودش نقطه مرکزی نیست.

- نقطه نویز (Noise Point): نقطه‌ای که نه مرکزی است و نه مرزی.

الگوریتم از یک نقطه مرکزی شروع می‌کند و با گسترش منطقه  $\epsilon$ ، تمام نقاط مرتبط را به کلاستر اضافه می‌کند. این فرآیند ادامه می‌یابد تا تمام نقاط مرکزی و نقاط مرتبط پوشش داده شوند.

نیازی به تعیین تعداد کلاسترها از قبل ندارد. می‌تواند کلاسترهایی با شکل‌های نامنظم را شناسایی کند. در برابر نویز مقاوم است.

انتخاب مناسب پارامترهای  $\epsilon$  و  $minPts$  دشوار است. در داده‌های با ابعاد بالا، عملکرد آن ممکن است کاهش یابد.

DBSCAN در دسته الگوریتم‌های کلاستریگ مبتنی بر چگالی قرار می‌گیرد. این دسته از الگوریتم‌ها کلاسترها را به‌عنوان مناطقی با چگالی بالا (نسبت به نقاط مجاور) تعریف می‌کنند و برای داده‌هایی که کلاسترهایی با شکل‌های پیچیده دارند، مناسب هستند.

## تفاوت DBSCAN و OPTICS

الف) DBSCAN بر اساس پارامترهای ثابت  $\epsilon$  و  $minPts$  عمل می‌کند. این محدودیت می‌تواند باعث شود که در داده‌هایی با چگالی‌های متغیر، برخی کلاسترها شناسایی نشوند. OPTICS نسخه پیشرفته‌تر DBSCAN است که به جای استفاده از یک مقدار ثابت  $\epsilon$ ، از یک رویکرد تطبیقی برای شناسایی کلاسترها استفاده می‌کند. این الگوریتم می‌تواند چگالی‌های متغیر را مدیریت کند.

ب) DBSCAN تنها یک مجموعه از کلاسترها را ایجاد می‌کند و نقاط نویز را مشخص می‌سازد. OPTICS به جای تولید مستقیم کلاسترها، یک ساختار مرتب‌شده از داده‌ها بر اساس چگالی تولید می‌کند. این ساختار می‌تواند برای استخراج کلاسترهای مختلف با تنظیمات متفاوت استفاده شود.

ج) DBSCAN اگر داده‌ها دارای کلاسترهایی با چگالی‌های متفاوت باشند، ممکن است برخی کلاسترها را از دست بدهد. OPTICS به دلیل استفاده از یک رویکرد انعطاف‌پذیر برای تعیین شعاع  $\epsilon$ ، می‌تواند کلاسترهایی با چگالی‌های متغیر را به خوبی شناسایی کند.

د) DBSCAN کلاسترهای مشخصی تولید می‌کند. OPTICS یک ترتیب (ordering) از نقاط تولید می‌کند که برای استخراج کلاسترهای مختلف استفاده می‌شود.