



موضوع هفته:

Clustering (K-Means, DBSCAN, Hierarchical) و معیارهای ارزیابی خوشبندی

سؤالات تحلیلی / ریاضی کوتاه

دو تمرین برنامه‌نویسی

بخش اول – سوالات تحلیلی / مفهومی

سؤال ۱

تابع هدف الگوریتم K-Means معمولاً به صورت مجموع مربعات فاصله‌ی نقاط هر خوشه از مرکز آن خوشه (WCSS) تعریف می‌شود. فرمول محاسبه مرکز جدید را از روی این تابع هدف به دست بیاورید.

سؤال ۲

روش Elbow یکی از روش‌های رایج برای انتخاب تعداد خوشه‌ها در K-Means است. توضیح دهید:

- در این روش چه مقداری بر حسب K روی نمودار رسم می‌شود و شکل کلی نمودار چه انتظاری داریم باشد؟
- منظور از " نقطه‌ی آرنج " چیست و چگونه می‌توان آن را به صورت شهودی تشخیص داد؟
- در چه شرایطی ممکن است نمودار WCSS آرنج واضحی نداشته باشد و این موضوع چه چیزی درباره‌ی ساختار داده به ما می‌گوید؟

سؤال ۳

در خوشبندی بدون ناظر، برچسبهای خوش (۰،۱،۰،۰...) هیچ معنای ذاتی خاصی ندارند و ممکن است بین اجرای الگوریتم یا نسبت به برچسبهای واقعی کلاس‌ها جابه‌جا شوند. توضیح دهید:

- چرا برای مقایسه خوششای K-Means با برچسبهای واقعی کلاس‌ها، معمولاً از الگوریتم Hungarian (یا روش‌های مشابه) استفاده می‌شود؟
- اگر بدون استفاده از Hungarian، خوششها را به صورت مستقیم با کلاس‌ها مقایسه کنیم، چه نوع خطای مفهومی یا عددی ممکن است رخ دهد؟
- نقش این الگوریتم در بهدست آوردن "بهترین نگاشت ممکن بین خوششها و کلاس‌ها" چیست؟

سؤال ۴

دو معیار NMI و Adjusted Rand Index (ARI) از معیارهای مهم برای ارزیابی کیفیت خوشبندی در مقایسه با برچسبهای واقعی هستند. مقایسه کنید:

- هر کدام به طور شهودی چه چیزی را اندازه‌گیری می‌کند؟
- هر دو معیار چگونه نسبت به جابه‌جایی برچسب خوششها (label permutation) ناپایا یا پایا هستند؟
- در چه شرایطی ممکن است مقدار ARI و NMI برای یک خوشبندی اختلاف قابل توجهی داشته باشد و این اختلاف چگونه قابل تفسیر است؟

سؤال ۵

معیارهای ارزیابی مانند ARI و NMI با دقت Accuracy (چه تفاوتی دارند؟ توضیح دهید:

- چرا استفاده‌ی مستقیم از دقت برای ارزیابی خروجی یک الگوریتم خوشبندی (بدون ناظر) چندان مناسب نیست؟
- چه ویژگی‌هایی در ARI و NMI وجود دارد که آن‌ها را برای سنجش کیفیت خوشبندی جذاب‌تر می‌کند (مثلاً تصحیح نسبت به شناسی یا در نظر گرفتن اطلاعات مشترک)؟
- اگر تنها به دقت تکیه کنیم، چه نوع سوءبرداشت‌هایی درباره خوب بودن یا بد بودن خوشبندی ممکن است به وجود آید؟

سؤال ۶

به صورت مفهومی تفاوت‌های اصلی بین سه روش DBSCAN، K-Means و Hierarchical Clustering خود حتماً به موارد زیر اشاره کنید:

- فرض هر الگوریتم درباره شکل خوش‌ها (کروی، چگالتی، سلسه‌مراتبی و ...)
- نیاز یا عدم نیاز به دانستن تعداد خوش‌ها قبل از اجرا
- حساسیت نسبت به مقیاس داده و پارامترها مثلًا `min_samples` در DBSCAN یا نوع `linkage` در Hierarchical
- رفتار هر الگوریتم در مواجهه با نویز و نقاط دورافتاده (`outlier`)

سؤال ۷

تابع هدف K-Means یک تابع غیرمحدب (non-convex) نسبت به تمام مراکز خوش به صورت همزمان است. توضیح دهید:

- این ویژگی چه تأثیری روی فرایند بهینه‌سازی و جواب نهایی الگوریتم دارد؟
- چرا اجرا کردن الگوریتم با نقاط اولیه مختلف initialization های متفاوت می‌تواند به جواب‌های متفاوت منجر شود؟
- در عمل چه راهکارهایی برای کاهش اثر مینیمم‌های محلی به کار می‌رود؟

بخش دوم – تمرین‌های برنامه‌نویسی (Python / sklearn)

نکته‌ی مهم: در هر دو تمرین این بخش، لازم نیست الگوریتم‌ها را از صفر پیاده‌سازی کنید؛ می‌توانید از توابع موجود در کتابخانه‌ی `sklearn` و سایر کتابخانه‌های مرسوم استفاده کنید. تمرکز باید بر طراحی آزمایش، تحلیل نتایج، انتخاب پارامترها و استفاده از معیارهای ارزیابی باشد، نه صرفاً نوشتن چند خط کد.

تمرین ۱ – خوشبندی داده‌ی تصاویر ارقام و ارزیابی با معیارهای پیشرفته

در این تمرین از دیتاست ارقام دست‌نویس (`digits`) موجود در کتابخانه `sklearn` استفاده کنید. این دیتاست شامل تصاویر 8×8 از ارقام ۰ تا ۹ است و به طور طبیعی ۱۰ کلاس دارد.

(الف) داده را بارگذاری کنید، شکل کلی آن (تعداد نمونه‌ها، تعداد ویژگی‌ها، تعداد کلاس‌ها) را گزارش دهید و در صورت نیاز مقادیر ویژگی‌ها را نرمال‌سازی یا استاندارد کنید. چند نمونه‌ی تصویری را به صورت شبکه‌ای نمایش دهید تا شهودی از داده به دست آورید.

ب) الگوریتم K-Means را برای چند مقدار متفاوت از K (برای مثال در بازه‌ای حول ۱۰، مثل ۶ تا ۱۴) روی این داده اجرا کنید.
برای هر مقدار K ، مقادیر زیر را محاسبه و ذخیره کنید:

- مقدار WCSS یا همان inertia گزارش شده توسط K-Means
- یک معیار داخلی مانند silhouette score برای خوشبندی بدون استفاده از برچسب‌های واقعی
- نمودار تغییرات این مقادیر را بر حسب K رسم کنید و از دید خودتان تحلیل کنید که آیا می‌توان " نقطه‌ی مناسب " برای K را از روی این نمودارها حدس زد یا خیر.

ج) یکی از مقدارهای K را که منطقی‌تر به نظر می‌رسد (مثلًاً نزدیک به ۱۰ و با رفتار مناسب در نمودارهای بخش قبل) انتخاب کنید
و K-Means را با این K چند بار و با initialization‌های متفاوت اجرا کنید. نشان دهید آیا مقدار تابع هدف و ساختار خوشبندی در اجراهای مختلف کاملاً یکسان است یا خیر و این رفتار را با بحث مینیمم‌های محلی مرتبط کنید.

د) با استفاده از برچسب‌های واقعی ارقام، برای خروجی خوشبندی K-Means برای همان K انتخابی معیارهای ارزیابی زیر را محاسبه کنید:

- Adjusted Rand Index (ARI)
- Normalized Mutual Information (NMI)

برای این کار ابتدا باید بین خوشبندی و برچسب‌های واقعی یک نگاشت مناسب برقرار کنید می‌توانید از الگوریتم Hungarian یا پیاده‌سازی آمده‌ی آن استفاده کنید. گزارش کنید پس از اعمال این نگاشت، دقت ظاهری (Accuracy) چقدر می‌شود و آن را در کنار ARI و NMI تفسیر کنید؛ توضیح دهید هر کدام از این اعداد چه زاویه‌ای از کیفیت خوشبندی را نشان می‌دهند و آیا بالابودن یکی از آن‌ها لزوماً به معنای بالابودن دیگری است یا خیر.

ه) در یک جمع‌بندی کوتاه، تحلیل کنید که K-Means روی این دیتاست چه نقاط قوت و چه محدودیت‌هایی دارد؛ به خصوص به این نکته اشاره کنید که بعضی ارقام از نظر شکل و شدت پیکسل‌ها شبیه هم هستند و این موضوع چگونه می‌تواند باعث تداخل خوشبندی شود، حتی اگر مدل به صورت عددی عملکرد نسبتاً خوبی داشته باشد.

تمرین ۲ - مقایسه‌ی Hierarchical Clustering و DBSCAN روی داده‌ی غیرکروی با نویز

در این تمرین می‌خواهید تفاوت رفتار دو الگوریتم DBSCAN و Hierarchical Clustering را روی داده‌های با شکل‌های غیرکروی و وجود نویز بررسی کنید. برای این کار می‌توانید از یکی از توابع تولید داده در `sklearn` مانند داده‌های دو هلالی معروف که توسط تابعی مثل `make_moons` تولید می‌شوند استفاده کنید و در صورت تمایل مقداری نویز تصادفی به داده اضافه نمایید.

(الف) یک دیتابست دوبعدی با حداقل دو خوش‌های غیرکروی (مثلًاً دو هلال در کنار هم) و تعداد مناسبی نقاط نویزی تولید کنید. داده را رسم کنید و توضیح دهید ساختار شهودی خوش‌های نویزها و نویزها چگونه است. در صورت نیاز داده را مقیاس‌بندی (scaling) کنید و دلیل این انتخاب را بیان کنید.

(ب) الگوریتم DBSCAN را با چند مجموعه پارامتر متفاوت مقادیر مختلف `min_samples` و `epsilon` روی این داده اجرا کنید. برای هر تنظیم، تعداد خوش‌های شناسایی شده، تعداد نقاط نویزی و شکل ظاهری خوش‌های را گزارش و روی نمودار نمایش دهید. نتایج را مقایسه کنید و توضیح دهید چگونه تغییر `min_samples` و `epsilon` می‌تواند به "بیش‌خوش‌بندی"، "کم‌خوش‌بندی" یا علامت‌گذاری بیش‌ازحد نقاط به عنوان نویز منجر شود.

(ج) روی همان داده، الگوریتم Agglomerative Hierarchical Clustering را با حداقل دو نوع `linkage` متفاوت مثلاً `single` و `ward` یا `complete` اجرا کنید. برای هر حالت، نمودار سلسله‌مراتبی (dendrogram) را رسم کنید و نشان دهید اگر در ارتفاع‌های مختلف آن را "برش" دهید، چه تعداد خوش و با چه ساختاری به دست می‌آید. توضیح دهید آیا این روش بدون دانستن تعداد خوش‌های هم می‌تواند به شما در انتخاب تعداد خوش‌های مناسب کمک کند یا خیر.

(د) نتایج Hierarchical Clustering و DBSCAN را از دو جنبه‌ی زیر با هم مقایسه کنید:

- توانایی در شناسایی شکل‌های غیرکروی و جدا کردن خوش‌های هلالی از یکدیگر
- توانایی در شناسایی و برچسب‌زندن نقاط نویزی
- در تحلیل خود توضیح دهید هر الگوریتم بر چه فرض‌هایی تکیه می‌کند و این فرض‌ها چرا در این نوع داده می‌توانند به تفاوت محسوس در خروجی‌ها منجر شوند.

(ه) در پایان، یک سناریوی واقعی (مثلًاً در حوزه‌ی تشخیص الگو، سامانه‌های فضایی، داده‌های سنسوری یا شبکه‌های اجتماعی) پیشنهاد کنید که در آن ترجیح دهید از DBSCAN استفاده کنید و یک سناریوی دیگر که در آن Hierarchical Clustering مناسب‌تر است. برای هر سناریو کوتاه توضیح دهید ساختار داده و نیاز مسئله چگونه با نقاط قوت آن الگوریتم هم‌خوانی دارد.