



## سوال ۱: معیارهای شباهت (تئوری)

همانطور که از مطالب درس آموختید، برای خوشه‌بندی داده‌ها نیاز به معیاری برای سنجش میزان شباهت بین آن‌ها داریم. تعیین این معیار شباهت به نوع داده‌ها و هدف مسئله بستگی دارد. فاصله اقلیدسی<sup>۱</sup> مثالی از یک معیار شباهت است که معمولاً برای داده‌های عددی استفاده می‌شود و شما با آن در طول درس آشنا شدید. در این سوال به دنبال آشنایی با برخی معیارهای شباهت دیگر و تمرین آموخته‌هایمان هستیم.

### قسمت الف: انتخاب معیار شباهت مناسب

در ستون سمت راست جدول ۱-۱ تعریف چند مسئله و در ستون سمت چپ چند معیار شباهت را مشاهده می‌کنید. هر مسئله را به معیار شباهت مناسب متصل کنید و دلیل انتخاب خود را شرح دهید.

معیار شباهت	مسئله
نقاط چگالی در دسترس <sup>۲</sup> (الگوریتم DBSCAN را مطالعه کنید).	خوشه‌بندی ستاره‌های آسمان
فاصله اقلیدسی	خوشه‌بندی اسناد متنی
شباهت جاکارد <sup>۳</sup>	خوشه‌بندی نتایج آزمایشات پزشکی به صورت اعداد باینری و بازه‌ی عددی
شباهت کسینوسی <sup>۴</sup>	خوشه‌بندی خانه‌های برای تعیین محل احداث مراکز پستی با وجود موانع

جدول ۱-۱: انتخاب معیار شباهت

### قسمت ب: ماتریس بی‌شباهتی<sup>۵</sup>

ماتریس بی‌شباهتی همه اطلاعات مربوط به تفاوت داده‌ها را، بر اساس یک معیار تفاوت  $d(.,.)$ ، در خود ذخیره می‌کند. برای  $n$  داده ماتریس بی‌شباهتی  $D$  عبارت است از:

$$D = (d_{ij}); d_{ij} = d(i, j), 1 \leq i, j \leq n$$

$$0 \leq d_{ij} \leq 1$$

بنابراین برای شباهت خواهیم داشت:

$$s_{ij} = 1 - d_{ij}$$

<sup>1</sup> Euclidian distance (L2 norm)

<sup>2</sup> Density reachable

<sup>3</sup> Jaccard similarity

<sup>4</sup> Cosine similarity

<sup>5</sup> Dissimilarity matrix

اکنون جدول ۱-۲ را در نظر بگیرید. داده‌های این جدول ترکیبی از ویژگی‌های عددی، اسمی و ترتیبی دارند. ماتریس بی‌شباهتی را برای این داده‌ها بدست آورید.

شماره	ویژگی ۱ (اسمی)	ویژگی ۲ (ترتیبی)	ویژگی ۳ (عددی)
۱	A	عالی	۴۵
۲	B	متوسط	۲۲
۳	C	خوب	۶۴
۴	A	عالی	۲۸

جدول ۱-۲

**نکته ۱:** برای معیار شباهت می‌توانید هم از مطالبی که در درس آموخته‌اید استفاده کنید و هم روش‌های دیگری به کار ببرید اما باید بتوانید روش خودتان را در گزارش توجیه کنید.

**نکته ۲:** برای ویژگی ترتیبی می‌توانید ابتدا با یک نگاشت، هر مقدار را به یک عدد تبدیل کنید و سپس از فاصله اقلیدسی استفاده کنید.

**نکته ۳:** در نهایت فاصله هر دو داده عبارت از میانگین فاصله ویژگی‌هایشان خواهد بود.

**سوال ۲: الگوریتم‌های خوشه‌بندی (تئوری)****قسمت الف: خوشه‌بندی با روش کا-میانگین<sup>۶</sup>**

با استفاده از روش خوشه‌بندی کا-میانگین، داده‌های جدول ۱-۲ را به ۲ خوشه تقسیم کنید. مراکز اولیه خوشه‌ها را  $B$  و  $C$  در نظر بگیرید. برای درک بهتری از مسئله، نقاط را رسم کرده و خوشه‌ها را روی شکل مشخص کنید. همچنین مشخص کنید نقطه  $x = (3, 2.5)$  در کدام دسته قرار می‌گیرد. برای فاصله از فاصله اقلیدسی استفاده کنید.

$i$	$x_1$	$x_2$
<b>A</b>	1	1
<b>B</b>	2	1
<b>C</b>	2	3
<b>D</b>	3	2
<b>E</b>	4	3
<b>F</b>	5	5

جدول ۱-۲: داده‌های مربوط به سوال کا-میانگین

**قسمت ب: خوشه‌بندی با روش سلسله مراتبی<sup>۷</sup>**

یکبار با استفاده از پیوند واحد<sup>۸</sup> و یکبار هم با پیوند کامل<sup>۹</sup>، با محاسبه فاصله اقلیدسی، داده‌های زیر را خوشه‌بندی کرده و نمودار درختی آن را رسم کنید.

$i$	$x_1$	$x_2$
<b>A</b>	0.45	<b>0.3</b>
<b>B</b>	0.22	<b>0.38</b>
<b>C</b>	0.08	<b>0.41</b>
<b>D</b>	0.26	<b>0.19</b>
<b>E</b>	0.35	<b>0.32</b>

جدول ۲-۲: داده‌های مربوط به سوال پیوند واحد

<sup>۶</sup> K-means clustering

<sup>۷</sup> Hierarchical clustering

<sup>۸</sup> Single-linkage

<sup>۹</sup> Complete-linkage

**سوال ۳: الگوریتم‌های خوشه‌بندی (پیاده‌سازی)****قسمت اول: الگوریتم کا-میانگین با انتخاب اولیه هوشمندانه**

در الگوریتم کا-میانگین ابتدا باید به طور تصادفی تعدادی مرکز خوشه انتخاب می‌کردیم. در درس آموختید که انتخاب مراکز خوشه اولیه ممکن است تاثیر زیادی روی نتایج نهایی الگوریتم و سرعت همگرایی آن داشته باشد. در این قسمت می‌خواهیم به پیاده‌سازی الگوریتم کا-میانگین با استفاده از یکی از روش‌های هوشمندانه برای انتخاب مراکز خوشه اولیه بپردازیم.

**الف)** الگوریتم کا-میانگین را با معیار فاصله اقلیدسی مطابق آنچه که در درس آموختید برای دادگان iris پیاده‌سازی کنید اما با این تفاوت که به جای انتخاب مراکز خوشه در ابتدا به صورت کاملاً تصادفی، از الگوریتم زیر استفاده کنید.

**الگوریتم جدید انتخاب مراکز خوشه اولیه:** اولین مرکز خوشه را به طور تصادفی و با احتمال یکنواخت از بین کل نقاط مجموعه داده‌ها انتخاب کنید. مراکز خوشه بعدی را باز هم به تصادف اما اینبار با احتمال متفاوتی انتخاب کنید. احتمال انتخاب هر مرکز خوشه (غیر از مرکز خوشه اول) متناسب با فاصله‌اش از نزدیکترین مرکز خوشه (نقطه‌ای که تا آن لحظه به عنوان مرکز خوشه انتخاب شده‌اند) است. برای فاصله هم می‌توانید از همان فاصله اقلیدسی استفاده کنید.

**ب)** الگوریتم پیاده‌سازی شده را برای تعداد خوشه‌های مختلف ۱ تا ۵ خوشه اجرا کنید. برای هر تعداد خوشه الگوریتم را به اندازه کافی تکرار کنید تا همگرایی حاصل شود. از مجذور میانگین فاصله اقلیدسی نقاط با مرکز خوشه‌هایشان به عنوان تابع هزینه استفاده کنید.

**نکته:** کل فرآیند خوشه‌بندی را چند بار با انتخاب مراکز خوشه اولیه متفاوت تکرار کنید و بهترین حالت را بر اساس مقدار تابع هزینه در نظر بگیرید.

**ج)** برای هر تعداد خوشه، نمودار تابع هزینه را بر حسب تکرارهای الگوریتم، تا رسیدن به همگرایی رسم کنید.

**د)** نمودار مقدار همگرایی تابع هزینه بر حسب تعداد خوشه‌ها را رسم کنید و سپس بر اساس آموخته‌هایتان در درس، بگویید که چه تعداد خوشه برای این دادگان مناسب است.

**تذکره ۱:** در این قسمت مجاز به استفاده از کتابخانه نیستید.

**تذکره ۲:** دادگان iris را می‌توانید از کتابخانه sklearn دریافت کنید.

**قسمت دوم: کا-میان ۱۰**

در روش کا-میان به جای فاصله اقلیدسی، از فاصله منهتن<sup>۱۱</sup> در تابع هزینه خود استفاده می‌کند. در این قسمت می‌خواهیم مقاومت این روش را با روش کا-میانگین نسبت به داده‌های تقلبی<sup>۱۲</sup> بررسی کنیم.

**الف)** قسمت قبل را برای روش کا-میان هم انجام دهید. (فقط به ازای ۳ خوشه)

<sup>10</sup> K-median

<sup>11</sup> Manhattan

<sup>12</sup> Outlier

ب) مراکز خوشه‌ها را که در قسمت‌های قبل با دو روش کا-میانگین و کا-میانه به ازای ۳ خوشه بدست آورده‌اید در نظر بگیرید. با توجه به ساختار دادگان iris و اندازه عددی ویژگی‌های هر داده، ۲۰ داده پرت به دادگان اضافه کنید. اکنون مجدداً فرآیندهای خوشه‌بندی کا-میانه و کا-میانگین را با حضور داده‌های تقلبی اجرا کنید. مقدار اختلاف مراکز خوشه را در حالتی که داده تقلبی نداشتیم با حالتی که داده تقلبی داریم، برای هر دو الگوریتم مقایسه کنید. کدام یک به داده‌های تقلبی حساس‌تر است؟

**نکات تحویل**

- ۱- مهلت تحویل بخش تئوری این تمرین **۱۵ آذرماه** و بخش پیاده‌سازی **۲۱ آذرماه** می‌باشد. (با توجه به اینکه بعد از تمرین پاسخ سوالات تئوری بارگذاری می‌شود امکان استفاده از گریس برای بخش تئوری را ندارید).
- ۲- انجام این تمرین به صورت یک‌نفره است.
- ۳- برای انجام این تمرین تنها مجاز به استفاده از زبان برنامه نویسی پایتون هستید.
- ۴- در صورت وجود تقلب نمره تمامی افراد شرکت کننده در آن **نمره صفر** لحاظ می‌شود.
- ۵- در صورتی که از منبعی برای هر بخش استفاده می‌شود، حتماً لینک مربوط به آن در گزارش آورده شود. وجود شباهت بین منبع و پیاده‌سازی در صورت ذکر منبع بلامانع است. اما در صورت مشاهده شباهت با مطالب موجود در سایت‌های مرتبط نمره کسر می‌گردد.
- ۶- نتایج و تحلیل‌های شما در روند نمره‌دهی دستیاران آموزشی تأثیرگذار است.
- ۷- لطفاً پاسخ تمرین خود را (به همراه کد/گزارش سوال کامپیوتری) به صورت زیر در صفحه درس آپلود نمایید:

HW[HW number]\_[Last\_name]\_[Student number].zip

- ۸- در صورت وجود هر گونه ابهام یا مشکل می‌توانید از طریق ایمیل با طراحان تمرین در تماس باشید:

- سیدمحمدمتین آل محمد: sm.matin.alemohammad@gmail.com
- نیلوفر فریدنی: nilu.1380@gmail.com