

به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



درس سیستم‌های هوشمند

تمرین شماره 3

آذر 1401

فهرست سوالات

- سوال ۱: تحلیلی 3
- الف) خوشه‌بندی به روش کا-میانگین 3
- ب) خوشه‌بندی سلسله‌مراتبی 4
- ب.۱) پیوند واحد 4
- ب.۲) پیوند کامل 4
- ب.۳) مقایسه 4
- سوال ۲: پیاده سازی الگوریتم خوشه بندی 5
- قسمت اول: خوشه‌بندی کا-میانگین ساده 5
- الف.۱) تاثیر تکرار خوشه‌ها 5
- الف.۲) تاثیر تکرار آزمایش 5
- قسمت دوم: خوشه بندی کا-میانگین هوشمند 6
- الف) طراحی الگوریتم 6
- ب) پیاده سازی الگوریتم 7
- ج) ارزیابی الگوریتم 7
- نکات تحویل: 8

سوال ۱: تحلیلی

در این بخش به بررسی و حل دستی الگوریتم‌های خوشه‌بندی^۱ می‌پردازیم.

الف) خوشه‌بندی به روش کا-میانگین^۲

با استفاده از روش خوشه‌بندی کا-میانگین ابتدا ماتریس فواصل (مربعات) داده‌های زیر را بدست آورید و سپس آنها را به دو خوشه تقسیم کنید. در نهایت تعیین کنید نقطه‌ی $x^* = (1,1)$ در کدام خوشه قرار می‌گیرد.

کا-روش به بندی خوشه های داده : ۱-۱ جدول
میانگین

i	x_1	x_2
A	0	0
B	0	1
C	-1	2
D	2	0
E	3	0
F	4	-1

نکته ۱: فاصله‌ی بین دو نقطه‌ی x_i و x_j بصورت زیر بدست می‌آید:

$$D_{ij} = \text{dist}_{\text{eucl.}}(x_i, x_j)^2$$

نکته ۲: نقاط ابتدایی و انتهایی داده‌ها را به عنوان مراکز اولیه خوشه‌بندی استفاده کنید.

نکته ۳: برای دستیابی به دید بهتری از مساله، نقاط را رسم کنید و مرکز خوشه‌ها را مشخص کنید.

¹ Clustering
² K-Means

ب) خوشه‌بندی سلسله‌مراتبی

جدول ۱-۲ ماتریس فواصل اقلیدسی نقاط را نشان می‌دهد.

ب.۱) پیوند واحد

با استفاده از روش پیوند واحد^۱، داده‌های زیر را خوشه‌بندی کرده و نمودار درختی آن را رسم کنید.

ب.۲) پیوند کامل

با استفاده از روش پیوند کامل^۲، داده‌های زیر را خوشه‌بندی کرده و نمودار درختی آن را رسم کنید.

ب.۳) مقایسه

دو مقدار از ماتریس داده‌شده را طوری تغییر دهید که نتیجه‌ی استفاده از پیوند واحد و پیوند کامل مشابه یکدیگر شوند. روند فکر خود را در گزارش بیاورید.

مراتبی‌سلسه روش به بندی‌خوشه فواصل ماتریس: ۱-۲ جدول

	A	B	C	D	E	F
A	0					
B	0.12	0				
C	0.51	0.25	0			
D	0.84	0.16	0.14	0		
E	0.28	0.77	0.7	0.45	0	
F	0.34	0.61	0.93	0.2	0.67	0

¹ Single Linkage
² Complete Linkage

سوال 2: پیاده سازی الگوریتم خوشه بندی

قسمت اول: خوشه بندی کا-میانگین ساده

در این سوال به پیاده سازی الگوریتم خوشه بندی کا-میانگین می پردازیم. هدف این الگوریتم بدون نظارت، کمینه کردن تابع هزینه^۱ است که به صورت مجموع فاصله ی نمونه های متعلق به هر خوشه تا مرکز آن خوشه (فاصله درون کلاسی) تعریف می شود.

با در نظر گرفتن مجموعه داده Iris سعی داریم الگوریتم خوشه بندی کا-میانگین ساده را با در نظر گرفتن نرم اقلیدسی پیاده سازی کنیم.

```
from sklearn.datasets import load_iris  
data = load_iris()
```

الف. ۱) تاثیر تکرار خوشه ها

پس از پیاده سازی الگوریتم، برای تعداد خوشه های ۵، ۱۰ و ۲۰، الگوریتم را ۲۰ بار تکرار کنید. نمودار مقدار هزینه را در حین اجرای الگوریتم رسم نمایید. آیا میتوانید تحلیل کنید که کدام تعداد خوشه مناسب تر است؟ (می توانید از معیار نسبت شباهت درونی^۲ به شباهت بیرونی^۳ استفاده نمایید).

الف. ۲) تاثیر تکرار آزمایش

برای تعداد خوشه های ۵، ۱۰ و ۲۰، الگوریتم را به تعداد کافی تکرار کنید. Ratio را به صورت رابطه ۱-۲، تعریف می کنیم.

$$Ratio = \frac{Outer\ distance}{Outer\ distance + Inner\ distance} \quad \text{رابطه ی ۱-۲.}$$

نمودار Ratio بر حسب تعداد دفعات تکرار الگوریتم را برای هر سه تعداد خوشه در یک نمودار رسم کنید. تحلیل کنید که کدام تعداد خوشه مناسب تر است؟

¹ Cost Function
² Intra-cluster Distance
³ Inter-cluster Distance

حال برای کل تعداد تکرارهای انجام شده به ازای تعداد خوشه‌های ۵، ۱۰ و ۲۰، میانگین و واریانس Ratio و تابع هزینه را محاسبه کنید. تحلیل کنید کدام تعداد خوشه با توجه به میانگین‌های به دست آمده، مناسب تر است. (برای درک بهتر می‌توانید، نمودار میانگین Ratio بر حسب تعداد خوشه و یا میانگین تابع هزینه بر حسب تعداد خوشه را رسم کنید، در نهایت با توجه به نمودار نتیجه‌گیری کنید).

قسمت دوم: خوشه بندی کا-میانگین هوشمند

در این قسمت، ناظری از بیرون اطلاعاتی در مورد قسمتی از دادگان^۱ به شما می‌دهد که شما را وادار می‌کند در فرآیند خوشه بندی محتاط تر عمل کنید. به این خوشه بندی، خوشه بندی شرطی^۲ گفته می‌شود. در خوشه بندی شرطی کا-میانگین تمامی مراحل مانند کا-میانگین ساده طی می‌شود با این تفاوت که در فرآیند تعلق هر داده به خوشه مربوطه نباید شرطی از شروط داده شده نقض گردد.

نگارش هر شرط داده شده به صورت سه تایی زیر می‌باشد:

$$[d_i, d_j, State] ; d_{i,j} \in D \text{ (Data set)}, State \in \{-1, 1\}$$

شرط بالا بیان می‌کند که هر دو داده دلخواه با اندیس^۳های متفاوت درون یک خوشه هستند یا خیر. در صورتی که عبارت State برابر 1 باشد، دو داده تحت هر شرایطی باید درون یک خوشه قرار بگیرند^۴ و در صورتی که برابر -1 باشد، دو داده تحت هر شرایطی باید در دو خوشه متفاوت قرار گیرند^۵. بدین صورت که یک داده به نزدیکترین خوشه‌ای تعلق می‌گیرد که هیچ کدام از دو نوع شرط ذکر شده نقض نشود.

الف) طراحی الگوریتم

الگوریتمی طراحی کنید که مرحله به مرحله مانند الگوریتم کا-میانگین ساده به سمت جواب بهینه همگرا شود.

¹ Background knowledge

² Constrained clustering

³ Index

⁴ Must-linked constraint

⁵ Cannot-link constraint

ب) پیاده سازی الگوریتم

در پوشه پروژه فایلی بنام "Constraint.txt" وجود دارد که 60 شروط برای پیاده سازی الگوریتم را در بر دارد. به ترتیب برای سه مقدار موجود در {20,40,60}، قسمتی از شروط رو به صورت رندوم انتخاب کنید و الگوریتم را اجرا کنید.

*** (اندیس‌های موجود در فایل قرار داده شده متناظر با اندیس‌های دادگان Iris خوانده شده توسط کتابخانه Scikit-learn می‌باشند)***

ج) ارزیابی الگوریتم

1. با توجه به معیارهای ارزیابی در درس، کیفیت خوشه بندی خود را برای سه حالت ذکر شده با رسم نمودار مقایسه و تحلیل کنید.
2. به طور خاص تر، برای تعداد 3 خوشه، به کمک برچسب موجود برای دادگان، دقت را برای دو حالت خوشه بند کا-میانگین و کا-میانگین هوشمند مقایسه کنید. آیا اطلاعات زمینه‌ای داده شده در حالت کا-میانگین هوشمند باعث افزایش دقت شده است ؟

نکات تحویل:

- مهلت تحویل این تمرین **17 آذر** می باشد.
- انجام این تمرین به صورت یک نفره است.
- برای انجام این تمرین تنها مجاز به استفاده از زبان برنامه نویسی پایتون هستید.
- در صورت وجود تقلب نمره تمامی افراد شرکت کننده در آن **-۱۰۰** لحاظ میشود.
- لطفا 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

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

مسئول تمرین سوال ۱: شیوا شاکری (shiva.shakeri@ut.ac.ir)

مسئول تمرین سوال ۲ قسمت اول: عاطفه ملاباقر (ati.mollabagher@ut.ac.ir)

مسئول تمرین سوال ۲ قسمت دوم: شایان واصف (sh.vassef@ut.ac.ir)