

۱- فرض کنید s یک معیار شباهت متريک تعريف شده روی مجموعه X باشد و داشته باشيم: $s(x,y) > 0, \forall x, y \in X$ چنانچه $f: R^+ \rightarrow R^+$ یک تابع پيوسته کاهشي يکنواخت باشد بطور يك داشته باشيم:

$$f(x) + f(y) \geq f\left(\frac{1}{\frac{1}{x} + \frac{1}{y}}\right), \quad \forall x, y \in R^+$$

نشان دهيد $f(s)$ یک معیار عدم شباهت متريک روی مجموعه X است.

۲- (اختياري) ثابت کنيد مقادير بيشينه و کمينه معیار شباهت فازي $s_F^q(\mathbf{x}, \mathbf{y}) = (\sum_{i=1}^l s(x_i, y_i)^q)^{1/q}$ بترتيب با $l^{1/q}$ و $l^{1/q}$ برابر است.

۳- (اختياري) شبه کد الگوريتم های خوشبندی ترتیبی (BSAS) و خوشبندی ترتیبی اصلاح شده (MBSAS) را با استفاده از یک معیار شباهت بجای معیار عدم شباهت بنویسید.

۴- (اختياري) الف- معادلات ۱۴.۵۶ و ۱۴.۵۷ کتاب تئودوريديس را برای الگوريتم پوسته های بیضیگونی فازی (FCES) بدست آورید. ب- روابط مربوط به تعیین پارامترهای الگوريتم FCES را بدست آورید.

۵- (اختياري) مجموعه داده زیر را در نظر بگيريد.

$$x_1 = (-1, 1), x_2 = (1, 1), x_3 = (-1, -1), x_4 = (1, -1)$$

می خواهیم از خوشبندی k-means با فاصله اقلیدسی استفاده کنیم با $K=2$. فرض کنید در ابتدا دو نقطه که به شکل تصادفی انتخاب می شوند به عنوان مرکز خوشبندی در نظر گرفته می شوند.

الف- تمام نتایج ممکن پس از اعمال الگوريتم خوشبندی را بنویسید.

ب- از بین همه نتایج ممکن، کدامیک کمترین هزینه را دارد؟ هزینه هر نتیجه را مجموع فاصله اقلیدسی هر نقطه از مرکز دسته ای که به آن نسبت داده شده است در نظر بگیرید.

ج- در شرایطی مثل حالت فوق، چطور می توان با استفاده از الگوريتم k-means به نتیجه های رسید که کمترین هزینه را داشته باشد؟

۶- (اختياري) مساله امتحان پایان ترم بازناسی آماری الگو سال ۱۳۸۹ در یک مساله خوشبندی از معیار مجموع مربعات خط استفاده می شود:

$$J = \sum_{i=1}^M J_i, \quad J_i = \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \hat{\mathbf{m}}_i\|^2$$

که در آن M تعداد خوشها و C_i مجموعه نمونهای خوش i -ام می‌باشد. چنانچه یک نمونه \mathbf{x} از خوش i -ام به خوش j -ام منتقل شود، نشان دهید مقادیر جدید مرکز ثقل خوشها و نیزتابع هزینه متضطر این خوشها به فرم زیر تغییر می‌یابند:

$$\begin{aligned}\hat{\mathbf{m}}'_j &= \hat{\mathbf{m}}_j + \frac{\mathbf{x} - \hat{\mathbf{m}}_j}{N_j + 1} & J'_j &= J_j + \frac{N_j}{N_j + 1} \|\mathbf{x} - \hat{\mathbf{m}}_j\|^2 \\ \hat{\mathbf{m}}'_i &= \hat{\mathbf{m}}_i - \frac{\mathbf{x} - \hat{\mathbf{m}}_i}{N_i - 1} & J'_i &= J_i - \frac{N_i}{N_i - 1} \|\mathbf{x} - \hat{\mathbf{m}}_i\|^2\end{aligned}$$

$\hat{\mathbf{m}}_j$ و $\hat{\mathbf{m}}_i$: مرکز ثقل خوشها قبل از انتقال نمونه \mathbf{x} از خوش i -ام به خوش j -ام
 $\hat{\mathbf{m}}'_j$ و $\hat{\mathbf{m}}'_i$: مرکز ثقل خوشها بعد از انتقال نمونه \mathbf{x} از خوش i -ام به خوش j -ام
 N_j و N_i : تعداد نمونهای خوشها j -ام و i -ام قبل از انتقال نمونه

۷- (اختیاری) مساله امتحان پایان ترم سال ۱۳۹۶

در فضای n بعدی فرض کنید N نمونه درون M خوش مجزا دسته‌بندی شده باشند و تابع هزینه خوشبندی که باید کمینه شود بصورت زیر تعریف شده باشد:

$$J_e = \sum_{i=1}^M \sum_{j=1}^{N_i} \|\mathbf{x}_{ij} - \mathbf{m}_i\|^2$$

که در آن M تعداد خوشها، N_i تعداد نمونهای خوش i -ام، \mathbf{m}_i مرکز ثقل نمونهای خوش i -ام و \mathbf{x}_{ij} نمونه j -ام از خوش i -ام می‌باشد.
 بازای $N > M$ نشان دهید که چنانچه یک خوش حاوی هیچ نمونه‌ای نباشد، این خوشبندی بهینه نیست. به عبارت دیگر برای حداقل شدن J_e ضروری است که حداقل یک نمونه در هر خوش وجود داشته باشد. (راهنمایی: ابتدا فرض کنید یکی از خوشها حاوی هیچ نمونه‌ای نیست و سپس نمونه‌ای را از یکی از خوشها به این خوش منتقل کنید.)

۸- (اختیاری) مساله امتحان پایان ترم ۱۳۹۷

در فضای n بعدی فرض کنید فاصله بین دو خوش A و B بصورت زیر تعریف شده باشد:

$$d(A, B) = \|\mathbf{m}_A - \mathbf{m}_B\|^2$$

که در آن \mathbf{m}_A و \mathbf{m}_B مرکز ثقل خوشها هستند. با استفاده از تعریف فوق چنانچه بخواهیم فاصله بین خوش K و یک خوش حاصل از ادغام دو خوش I و J را محاسبه کنیم نشان دهید که می‌توان نوشت:

$$d(K, I + J) = \frac{N_I}{N_I + N_J} d(I, K) + \frac{N_J}{N_I + N_J} d(J, K) - \frac{N_I N_J}{(N_I + N_J)^2} d(I, J)$$

که در آن N_I و N_J به ترتیب تعداد نمونهای خوش I و J هستند.

۹- (اختیاری) سوال امتحانی پایان ترم ۱۳۹۸

در فضای n بعدی فرض کنید N نمونه در اختیار است و می‌خواهیم با استفاده از الگوریتم K-means آنها را در K خوشبندی نماییم. مطابق این الگوریتم، مجموعه بردارهای مرکز خوشها $\{\mu_1, \mu_2, \dots, \mu_K\}$ با کمینه‌سازی فاصله میان نمونه‌ها از نزدیکترین مرکز خوش بدست می‌آید. به عبارت دیگر در این روش تابع هزینه زیر کمینه می‌شود:

$$L(\boldsymbol{\mu}) = \sum_{i=1}^N \min_{j \in \{1, \dots, K\}} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$$

و اختصاص هر نمونه \mathbf{x}_i به خوش j -ام بر اساس متغیر $Z_i = \arg \min_{j \in \{1, \dots, K\}} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$ انجام می‌شود. یعنی Z_i معرف برچسب نزدیکترین مرکز خوش به نمونه \mathbf{x}_i می‌باشد. الگوریتم K-means بصورت تکراری در دو مرحله‌ی بهروزسازی برچسب‌های Z_i (مرحله‌ی

برچسب زنی) و بهروزسازی مراکز خوشها ($\mu_j = \frac{1}{|\{i:z_i=j\}|} \sum_{i:z_i=j} x_i$ (مرحله بازنظمیم) اجرا می شود. الگوریتم هنگامی متوقف می شود که در مرحله‌ی برچسب زنی تغییری در برچسب نمونه‌ها بوجود نیاید. نشان دهید که الگوریتم K-means همواره به یک نقطه کمینه (محلی یا سراسری) همگرا می شود.

راهنمایی: کافی است نشان دهید کهتابع هزینه بصورت تکراری تا رسیدن به نقطه همگرایی بصورت یکنواخت کاهش می‌یابد. این روند کاهشی در هر تکرار الگوریتم را برای مرحله‌ی برچسب زنی و مرحله بازنظمیم بطور جداگانه نشان دهید.

۱۰- سوال امتحانی پایان ترم ۱۳۹۹

تابع هزینه برای مساله خوش بندی k -میانگین با k خوش، و نمونه‌های x_1, \dots, x_n و مراکز خوش μ_1, \dots, μ_k بصورت زیر تعریف شد:

$$L = \sum_{j=1}^k \sum_{x_i \in S_j} \|x_i - \mu_j\|^2$$

که در آن S_j زیرمجموعه نمونه‌هایی است که به μ_j نسبت به سایر مراکز خوش نزدیک‌تر هستند.

الف- به جای به روز سازی μ_j بر اساس متوسطگیری از نمونه‌های z_k می‌خواهیم از روش گرادیان نزولی دسته‌ای (batch gradient descent) که در آن زیرمجموعه‌های S_j ثابت نگه داشته می‌شوند استفاده کنیم. رابطه به روز سازی μ_j را چنانچه نرخ آموزش ϵ باشد بدست آورید.

ب- رابطه به روز سازی μ_j را بر اساس الگوریتم گرادیان نزولی برهنخ (به ازای یک تک نمونه) و نرخ آموزش ϵ بنویسید.

ج- در این قسمت می‌خواهیم معادل بودن رابطه بدست آمده در قسمت الف را با الگوریتم استاندارد k -میانگین بررسی کنیم. در الگوریتم استاندارد، میانگین نمونه‌های هر خوش را به عنوان مرکز آن خوش در نظر می‌گیریم. به نظر می‌رسد چنانچه در الگوریتم بدست آمده در بند الف، یک مقدار خاص برای نرخ آموزش (یعنی ϵ) در نظر گرفته شود، رابطه به روز سازی مراکز خوش ها در هر دو روش یکسان خواهد شد. چنانچه μ_1 مرکز خوش شماره ۱ باشد، مقدار ϵ را طوری تعیین کنید که رابطه به روز سازی مرکز این خوش در هر دو روش یکسان شود. (توجه شود که مقدار ϵ می‌تواند بطور خاص برای به روز سازی خوش شماره ۱ بدست آید و برای دیگر خوش ها متفاوت باشد).

۱۱- (اختیاری) فرض کنید داده‌هایی به ما داده شده که شامل چند کلاس مختلف است و هر کلاس دارای یک توزیع احتمال متفاوت می‌باشد. برچسب نمونه‌ها در دسترس نیست و قرار است از k-means به منظور خوشبندی استفاده شود. شرایطی که باعث تضعیف اثربخشی k-means می‌گردند را علامت بزنید.

الف. برخی کلاس‌ها دارای توزیع نرمال نباشند

ب. واریانس توزیع‌ها از تمام جهات دارای واریانس کمی باشد

ج. میانگین کلاس‌ها با هم برابر باشد.

د. انتخاب شود که n تعداد نمونه‌های در دسترس است.

۱۲- (پایان ترم ۱۴۰۲) معیار زیر برای خوشبندی بردارهای \mathbf{x} در داخل C خوش بکار می‌رود:

$$J_G = \sum_{i=1}^c \sum_{\mathbf{x} \in R_i} (\mathbf{x} - \mathbf{m}_i)^t \Sigma_T^{-1} (\mathbf{x} - \mathbf{m}_i)$$

در رابطه فوق:

$$\mathbf{m}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in R_i} \mathbf{x} \quad \Sigma_T = \sum_{\mathbf{x} \in X} (\mathbf{x} - \mathbf{m}) (\mathbf{x} - \mathbf{m})^t \quad \mathbf{m} = \frac{1}{N} \sum_{\mathbf{x} \in X} \mathbf{x}$$

در روابط فوق X معرف مجموعه کل نمونه‌ها، \mathbf{m} معرف مرکز ثقل عمومی، Σ_T ماتریس کواریانس عمومی، R_i ناحیه مربوط به خوشة i ام، N_i تعداد نمونه‌های قرار گرفته در خوشه i ام و \mathbf{m}_i مرکز ثقل خوشه i ام را مشخص می‌کند. حداقل‌سازی معیار فوق خوشه‌بندی بر اساس حداقل مربعات فاصله با معیار فاصله ماهالانوبیس را نتیجه می‌دهد. نشان دهید اگر یک نمونه \mathbf{y} از خوشه i به خوشه j منتقل شود، مقدار تابع هزینه بصورت زیر تغییر می‌کند:

$$J'_G = J_G + \left[\frac{N_j}{N_j + 1} (\mathbf{y} - \mathbf{m}_j)^t \Sigma_T^{-1} (\mathbf{y} - \mathbf{m}_j) - \frac{N_i}{N_i - 1} ((\mathbf{y} - \mathbf{m}_i)^t \Sigma_T^{-1} (\mathbf{y} - \mathbf{m}_i)) \right]$$

(۱۴۰۳) - پایان ترم

الف- نتایج خوشه‌بندی (اعضای هر خوشه و مراکز خوشه‌ها) با استفاده از الگوریتم k -means مجموعه داده زیر، پس از یک تکرار را گزارش دهید. (تعداد خوشه‌ها را ۳ در نظر بگیرید و در ابتدا داده‌های ۱، ۴ و ۷ را مراکز خوشه در نظر بگیرید و از فاصله اقلیدسی استفاده کنید).
 $A_1 = (2, 10)$, $A_2 = (2, 5)$, $A_3 = (8, 4)$, $A_4 = (5, 8)$, $A_5 = (7, 5)$, $A_6 = (6, 4)$, $A_7 = (1, 2)$, $A_8 = (4, 9)$.
در زیر ماتریس فاصله اقلیدسی داده‌ها نسبت به یکدیگر، جهت افزایش سرعت حل این مسئله، به شما داده شده است.

	A1	A2	A3	A4	A5	A6	A7	A8
A1	0	$\sqrt{25}$	$\sqrt{36}$	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{65}$	$\sqrt{5}$
A2		0	$\sqrt{37}$	$\sqrt{18}$	$\sqrt{25}$	$\sqrt{17}$	$\sqrt{10}$	$\sqrt{20}$
A3			0	$\sqrt{25}$	$\sqrt{2}$	$\sqrt{2}$	$\sqrt{53}$	$\sqrt{41}$
A4				0	$\sqrt{13}$	$\sqrt{17}$	$\sqrt{52}$	$\sqrt{2}$
A5					0	$\sqrt{2}$	$\sqrt{45}$	$\sqrt{25}$
A6						0	$\sqrt{29}$	$\sqrt{29}$
A7							0	$\sqrt{58}$
A8								0

ب- شبه کد انجام محاسبات بند الف را در حالت کلی یعنی C خوشه و تعداد نمونه‌های ورودی برابر N بنویسید و در آن ورودی‌ها، خروجی‌ها و حلقه عملیات تا رسیدن به نتیجه را مشخص کنید.

موفق باشد