

کلیف ۳ یا رنگی داشتن

مردی کاف ۹۹۲۱۰۷۵۳

۱.۱  
الف) اگر بخواهیم  $k$ -means حالت خاص از خوشه بندی با GMM است به گونه ای که پارامترها یکسان باشند و

کواریانسها صفر و mixture weightها هم برابر باشند. پس فرض  $k$ -means این است که تمام کلاسترها

هموز هستند و این در بسیاری از موارد پاسخ بهینه را نمیدهد. در Gaussian Mixture برای تخمین پارامترها معمولاً

از MLE استفاده میکنیم. برای داده های  $D$  و پارامترهای  $(\pi_k, \mu_k, \Sigma_k)$  از برای هر cluster از  $k$  تخمین بزنیم و از رابطه زیر استفاده میکنیم.

$$\log P(D | \pi, \mu, \Sigma) = \sum_i \log \sum_k \pi_k P_n(z_i = k) P_{\mu, \Sigma}(x_i | z_i = k)$$

در اینجا اگر پارامترها را به صورت  $\Sigma_k = \sigma_k^2 I$  در نظر بگیریم یعنی فرض کرده ایم که تمام کلاسترها به شکل دایره باشند و این EM با  $k$  خوشه را به  $k$ -means تبدیل کرده ایم.

۱.۱ ب) در خوشه بندی سلسله مراتبی با روش مجامعه شباهت به صورت Max یا Complete linkage ممکن است

که داده ای در یک خوشه داده ای در خوشه دیگر قرار بگیرد از داده ای در خوشه خودش باشد زیرا که این روش عامل داده کلاستهای

بزرگ را بشکند و همچون به outlierها اهمیت بسیاری میدهد و ممکن است این نوع خطا در

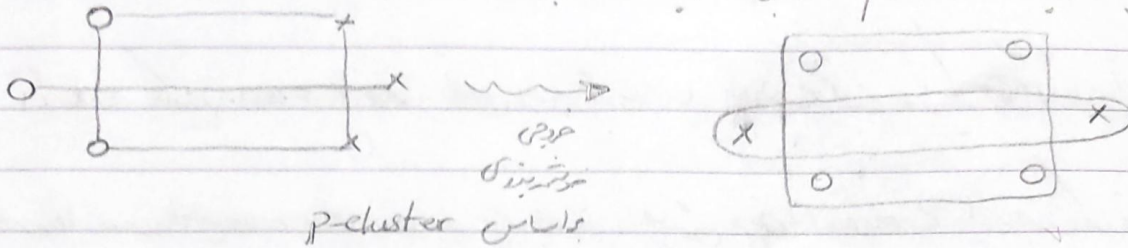


خوشه بندی نشان داده شده است:

۱.۱ ج. این نوع از خوشه بندی با توجه به توصیفات مشابه به نظر می رسد که بطل استفاده کردن از معیار فاصله نمی تواند باشد.

صورت خاص برود زیرا که اگر داده ها به صورت زیر باشند که x های ۱ ده و ۵ های ۱ ده به نظر می رسد که

بال خاص صورت زیر باشند با  $P=100$  هم، خاص مطلوب نخواهد بود.



۱.۲

الف) در صورت Single-linkage رابطه انعطاف پذیر و مستقیم باید هرگاه که شکل می دهند زیرا که بسیار به یکدیگر نزدیک اند.

سپس نقاط بین دسته های (۳ و ۴) و (۱ و ۲) باید در کلاس قرار بگیرند که با روش کوتاه ترین فاصله دو دسته ۳ و ۴ با

یکدیگر و دو دسته ۱ و ۲ نیز باید یک کلاس را تشکیل می دهند. کلاس های نهایی: (۳ و ۴) و (۱ و ۲)

ب) در صورت Complete-linkage رابطه انعطاف پذیر هر یک از دسته های ۱ تا ۴ باید یک خوشه می شوند و سپس در نهایت

به حالتی می رسیم که فاصله دورترین نقطه از خوشه ۱ تا دورترین نقطه از خوشه ۲ بیشترین فاصله بین خوشه ۱ و ۳ می شود.

بنابراین دو خوشه (۱ و ۳) و (۲ و ۴) ایجاد می شوند. کلاس های نهایی: (۲ و ۴) و (۱ و ۳)

ج) در صورت average-linkage سپس از اینک اعضای داخل کلاس ها باید یک شکل را داشته باشند و فاصله بین خوشه ها

بین (۱ و ۳) و (۲ و ۴) میانگین بین (۱ و ۲) خواهد شد. کلاس ها به صورت (۲ و ۴) و (۱ و ۳) می شوند.



۱.۲) در شکل (ب) معیار Single-link می‌تواند معقول عمل کند زیرا که در ابتدا نقاطی که به هم یک از منتهی می‌باشند

و به تدریج به تدریج با یکدیگر یک خوشه تشکیل داده و سپس این خوشه‌ها گسترش می‌یابند تا کل ترکیب از خوشه‌ها  
پوشانند و این صورت در بعضی موارد یک خوشه جابجایی می‌شوند.

در شکل (ج) به دلیل وجود نیزه‌ها در معیار Single-link خوب عمل نمی‌کند زیرا که ممکن است دو خوشه با یکدیگر

ترکیب کند و از این دو معیار Complete-link و average-link می‌تواند با معیار average-linkage با معیار بهتری نتایج حاصل کرد.

$$\left. \begin{array}{l} \mu = \text{mean}(C) \\ \mu' = \text{mean}(C') \\ \bar{\mu} = \text{mean}(C \cup C') \end{array} \right\} \Rightarrow \text{cost}(C \cup C') - \text{cost}(C) - \text{cost}(C')$$

$$= \sum_{x \in C \cup C'} \|x - \bar{\mu}\|^2 - \sum_{x \in C} \|x - \mu\|^2 - \sum_{x \in C'} \|x - \mu'\|^2$$

$$= \sum_{x \in C} (\|x - \bar{\mu}\|^2 - \|x - \mu\|^2) + \sum_{x \in C'} (\|x - \bar{\mu}\|^2 - \|x - \mu'\|^2)$$

$$= |C| \cdot \|\mu - \bar{\mu}\|^2 + |C'| \cdot \|\mu' - \bar{\mu}\|^2$$

$$\bar{\mu} = (|C|\mu + |C'|\mu') / (|C| + |C'|) \Rightarrow = |C| \cdot \left\| \frac{|C'|}{|C| + |C'|} (\mu' - \mu) \right\|^2 + |C'| \cdot \left\| \frac{|C|}{|C| + |C'|} (\mu - \mu') \right\|^2$$

$$= \frac{|C| \cdot |C'|}{|C| + |C'|} \|\mu - \mu'\|^2$$

	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	0.70	0.45	0	
F	0.34	0.61	0.93	0.20	0.67	0

Single-linkage  $\bar{C}$   $C1.F$

	A, B	C	D	E	F
A, B	0				
C	0.25	0			
D	0.16	✓0.14	0		
E	0.28	0.70	0.45	0	
F	0.34	0.93	0.20	0.67	0

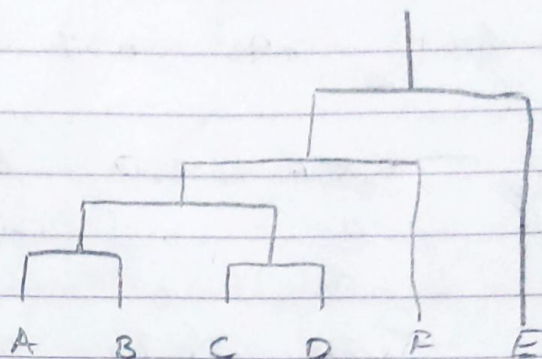
	A, B	C, D	E	F
A, B	0			
C, D	✓0.16	0		
E	0.28	0.45	0	
F	0.34	0.20	0.67	0

	A, B, C, D	E	F
A, B, C, D	0		
E	0.28	0	
F	✓0.20	0.67	0



	A, B, C, D, F	E
A, B, C, D, F	0	
E	10.28	0
	A, B, C, D, F, E	
A, B, C, D, F, E	0	

Dendrogram:



Complete linkage

	A	B	C	D	E	F
A	0					
B	10.12	0				
C	0.51	0.25	0			
D	0.84	0.16	0.14	0		
E	0.28	0.77	0.70	0.45	0	
F	0.34	0.61	0.93	0.20	0.67	0

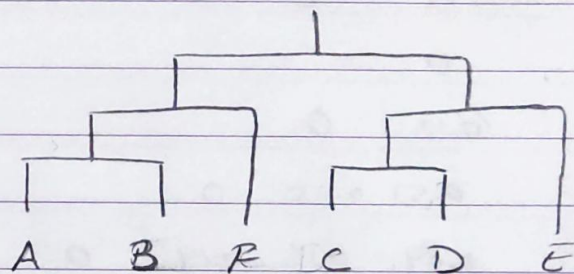
	A, B	C	D	E	F
A, B	0				
C	0.51	0			
D	0.84	10.14	0		
E	0.77	0.70	0.45	0	
F	0.61	0.93	0.20	0.67	0

	A, B	C, D	E	F
A, B	0			
C, D	0.84	0		
E	0.77	0.70	0	
F	0.61	0.93	0.67	0

	A, B, F	C, D	E
A, B, F	0		
C, D	0.93	0	
E	0.77	0.70	0

	A, B, F	C, D, E
A, B, F	0	
C, D, E	0.93	0

Dendrogram:



ج. با تغییر مقدار فاصله A, F از 0.34 به 0.15 و تغییر فاصله D, E از 0.45 به 0.15، نمودار dendrogram

حرفه‌ای فاصله نمودار نسبت به خواهد شد. در صورتی برای تغییر در فاصله فواصل، فاصله به نزدیک‌ترین محاسبه می‌شود.

استی ب شود به عنوان مقدار کمینه برای یکی کردن آن سطرهاست.



۲۰۱) TNoM یک روش Non-parametric است و در واقع به صورت بیشترین دفعات در داده‌ها از توزیع خاصی نمی‌دانند

روش کار آن بدین صورت است که فرض کنیم که تعدادی از داده‌ها را داریم و به آن label های

+ و - را می‌دهیم و اینصورت آن فرض کنیم که a با b تفاوت دارد label مثبت خواهد شد و با b تفاوت ندارد

label منفی خواهد شد و اینصورت مقدار R به اینصورت تعریف می‌شود که میزان بیان آن و در هر یک از داده‌ها به صورت مثبت و

منفی به نداد است  $\langle R_1, R_2, \dots, R_n \rangle$  که R نشان دهنده بیشترین بیان و است. سپس به هر rank

به اینصورت تعریف می‌شود که نشان دهنده جایگاه‌های هر داده اگر بیان متعلق به بافت کرده یا بود مقدار + و اگر متعلق به نبود مقدار -

- را قرار دهد. به طریقی که بیان و بافت داده به صورت  $\{10, 20, 30, 50, 80\}$  و جایگاه داده به صورت

$\{40, 70\}$  باشد به نظر رتبه ۷ به صورت بزرگ خواهد بود  $\{+, +, +, -, +, -, +\}$  و  $v = 0$  سپس

برای تمام پارتیشن‌های اول می‌توانیم  $v$  را بدست آوریم score محاسبه می‌شود که در نهایت کمینه این مقادیر

برای TNoM Score خواهد بود. حال score های اینصورت محاسبه می‌شوند که برای ۲ پارتیشن اول یک بار فرض می‌کنیم که آنها

کلاس + و کلاس منفی و بارها برعکس این فکر می‌کنیم و در هر بار مقدار داده‌های از خانواده دیگر در پارتیشن را می‌شماریم و این

دو عدد که برای انواع تقاطع ایجاد شود min می‌گیریم و سپس بین کلاص مقادیر پارتیشن‌ها min می‌گیریم که همان

TNoM Score می‌شود

$$TNoM(v) = \min_{x \neq y} \min ( [\#_-(x) + \#_+(y)], [\#_+(x) + \#_-(y)] )$$

محدود تعداد میں کڑوں کا بلور انتخاب کے مقدار Threshold نامی جان بوجھ کر حد میں عدد misclassification کا

اس کا Threshold Number of Misclassification نیز نشان لکھیں وضع 20

۱۲) با انتخاب ویژگی از بین ویژگی‌های مهم و رایج ما بعضی ویژگی‌ها را انتخاب می‌کنیم و بعضی دیگر را قایل  
تفسیر می‌کنیم و به طور مثال می‌توانیم بگوییم که هرگاه از ویژگی‌های مهمی که قبلاً انتخاب کرده‌ایم بخواهیم ویژگی‌های  
استخراج ویژگی‌ها را یکی از ویژگی‌های مهم و رایج و مهم که می‌توانیم بگوییم که این ویژگی‌ها را می‌توانیم به روش‌های  
انتخاب ویژگی‌ها، تفسیر کنیم و به روش‌های مختلف می‌توانیم به روش‌های استخراج ویژگی‌ها.

۲۰۲۰

Pearson correlation: عیب اصلی این روش این است که فقط احساس به روابط خطی است و آن ها را تشخیص میدهد

طریقه آزمون: این معیار را برای  $\chi^2$  و نظر کنیم هر چند که انطباق یک به یک برقرار است ولی P. correlation

مقایسه نزدیک به صفر برقرارند؟ بنابراین توصیه میشود که تنها به معیار P. correlation بسنده نکنیم و همواره

داریم و این را بسنجیم.

bin بندی هم به صورتی معیار بسیار واسطه bin بندی می شود.



سوال ۳)

۳.۳) ریشه روش  $V = WH$  که  $V$  بردارهای فضا است  $W$  بردارهای تبدیل یافته و ماتریس  $H$  ماتریس

خبر است؟ بنا بر این بردارهای تبدیل یافته ماتریس  $W$  طی فرایند