

Q1. Cosmetics Purchases

Trans. #	Bag	Blush	Nail Polish	Brushes	Concealer	Eyebrow Pencils	Bronzer
1	0	1	1	1	1	0	1
2	0	0	1	0	1	0	1
3	0	1	0	0	1	1	1
4	0	0	1	1	1	0	1
5	0	1	0	0	1	0	1
6	0	0	0	0	1	0	0
7	0	1	1	1	1	0	1
8	0	0	1	1	0	0	1
9	0	0	0	0	1	0	0
10	1	1	1	1	0	0	0
11	0	0	1	0	0	0	1
12	0	0	1	1	1	0	1

- a. The matrix form shows that each row represents a transaction, and the columns are different cosmetic products. The values of the matrix form show whether that kind of the cosmetics is in the transaction or not. For example: In transaction 6, the customer bought the concealer, so the value is 1. Because the customer didn't buy like bag, brushes, the values would be 0

TABLE 14.15 ASSOCIATION RULES FOR COSMETICS PURCHASES DATA

lhs	rhs	support	confidence	lift
1 {Blush, Concealer, Mascara, Eye.shadow, Lipstick} => {Eyebrow.Pencils}	0.013	0.3023255814	7.198228128	
2 {Trans., Blush, Concealer, Mascara, Eye.shadow, Lipstick} => {Eyebrow.Pencils}	0.013	0.3023255814	7.198228128	
3 {Blush, Concealer, Mascara, Lipstick} => {Eyebrow.Pencils}	0.013	0.2888888889	6.878306878	
4 {Trans., Blush, Concealer, Mascara, Lipstick} => {Eyebrow.Pencils}	0.013	0.2888888889	6.878306878	
5 {Blush, Concealer, Eye.shadow, Lipstick} => {Eyebrow.Pencils}	0.013	0.2826086957	6.728778468	
6 {Trans., Blush, Concealer, Eye.shadow, Lipstick} => {Eyebrow.Pencils}	0.013	0.2826086957	6.728778468	

b.

- i. The confidence of the first row means that given we see blush, concealer, mascara, eye shadow and lipstick in a transaction, what's the probability that we see eyebrow pencil? The calculation would be :

$$\text{Confidence} = P(\text{eyebrow pencil} | \text{blush, concealer, mascara, eyeshadow, lipstick})$$

$$= \frac{(\text{number of transactions with eyebrow pencil, blush, concealer, mascara, eye shadow, lipstick})}{(\text{number of transactions with blush, concealer, mascara, eye shadow, lipstick})}$$

$$= 0.3023$$

II. Support ensures that items (blush, concealer, mascara, eye shadow, lipstick) that occur relatively frequently in transactions. The calculation would be :

$$\text{Support} = P(\text{blush, concealer, mascara, eyeshadow, lipstick})$$

$$= \frac{\text{(number of transactions with blush, concealer, mascara, eye shadow, lipstick)}}{\text{(total number of transactions)}}$$

$$= 0.0130$$

III. Lift examines how much the rule improves occurrence of the consequent item without the rule. That is, the probability buying consequent increases x percent when we see antecedent. The calculation would be :

$$\text{Lift} = \frac{P(\text{eyebrow pencil} | \text{blush, concealer, mascara, eye shadow, lipstick})}{P(\text{eyebrow pencil})}$$

$$= \frac{P(\text{eyebrow pencil, blush, concealer, mascara, eye shadow, lipstick})}{P(\text{eyebrow pencil})} \times P(\text{blush, concealer, mascara, eye shadow, lipstick})$$

$$= 7.1982$$

IV. If blush, concealer, mascara, eye shadow and lipstick are purchased, then we are 30.23% confident that eyebrow pencil will also be purchased. This rule would be 719% better than purchasing eyebrow pencil only.

C.

- i. The first rule of the output shows that if brushes are purchased, then we assure that the nail polish will also be purchased. Also, It would be 357% better than purchasing nail polish only.

Look at the second and third rule, both the antecedents include blush and eye shadow, and the consequent is also the same(mascara). These two rules both show over 90% confidence. The only difference between these two is that the antecedent of the second rule adds concealer.

ii.

	lhs	rhs	support	confidence	lift	count
[1]	{Brushes}	=> {Nail.Polish}	0.149	1.0000000	3.571429	149
[2]	{Blush,Concealer,Eye.shadow}	=> {Mascara}	0.119	0.9596774	2.688172	119
[3]	{Blush,Eye.shadow}	=> {Mascara}	0.169	0.9285714	2.601040	169

[2],[3] 只差[3]多了Concealer

[4]	{Nail.Polish,Eye.shadow}	=> {Mascara}	0.119	0.9083969	2.544529	119
[5]	{Concealer,Eye.shadow}	=> {Mascara}	0.179	0.8905473	2.494530	179

[2],[5] 重複只差[2]多了Blush

[6]	{Bronzer,Eye.shadow}	=> {Mascara}	0.124	0.8794326	2.463397	124
[12]	{Bronzer,Mascara}	=> {Eye.shadow}	0.124	0.9051095	2.375615	124

[6],[12] 重複只是順序顛倒

[7]	{Concealer,Eye.shadow,Eyeliner}	=> {Mascara}	0.114	0.8769231	2.456367	114
[8]	{Blush,Mascara}	=> {Eye.shadow}	0.169	0.9184783	2.410704	169
[9]	{Eye.shadow,Lipstick}	=> {Mascara}	0.110	0.8527132	2.388552	110
[10]	{Mascara,Lipstick}	=> {Eye.shadow}	0.110	0.9090909	2.386065	110

[9],[10] 重複只是順序顛倒

[11] {Blush,Concealer,Mascara} => {Eye.shadow} 0.119 0.9083969 2.384244 119

刪除多餘的規則後如下：

lhs	rhs	support	confidence	lift	count
[1] {Brushes}	=> {Nail.Polish}	0.149	1.0000000	3.571429	149
[2] {Blush,Concealer,Eye.shadow}	=> {Mascara}	0.119	0.9596774	2.688172	119
[4] {Nail.Polish,Eye.shadow}	=> {Mascara}	0.119	0.9083969	2.544529	119
[6] {Bronzer,Eye.shadow}	=> {Mascara}	0.124	0.8794326	2.463397	124
[7] {Concealer,Eye.shadow,Eyeliner}	=> {Mascara}	0.114	0.8769231	2.456367	114
[9] {Eye.shadow,Lipstick}	=> {Mascara}	0.110	0.8527132	2.388552	110

這些規則可能是比較重要且沒有重複性的的，當然購物籃分析中，左邊所包含的商品及右邊推薦的商品，其相關性可能也是一件能夠參考的事情，如[2]與[3]規則中，

[2] {Blush,Concealer,Eye.shadow} => {Mascara}
[3] {Blush,Eye.shadow} => {Mascara}

[3] 其實也有其重要的涵義，有 {Blush,Eye.shadow} 但沒有 {Concealer,Eye.shadow}，可以說明 {Blush,Eye.shadow} 的關聯性較佳，並且通常買這兩個商品的人也會買 {Mascara}，而[9]跟[10]也有同樣概念，有

[9] {Eye.shadow,Lipstick} => {Mascara}
[10] {Mascara,Lipstick} => {Eye.shadow}

但卻沒有 {Mascara, Eye.shadow} => {Lipstick}，表示 {Mascara, Eye.shadow} 兩者的關聯性可能較不高。

Q2. Course ratings

TABLE 14.16 RATING OF ONLINE STATISTICS COURSES: 4 = BEST, 1 = WORST, BLANK = NOT TAKEN

	SQL	Spatial	PA 1	DM in R	Python	Forecast	R Prog	Hadoop	Regression
L N	4				3	2	4		2
M H	3	4			4				
J H	2	2							
E N	4			4			4		3
D U	4	4							
F L		4							
G L		4							
A H		3							
S A			4						
R W			2					4	
B A			4						
M G			4			4			
A F			4						
K G			3						
D S	4			2			4		

a.

The equation below is the correlation between two users' ratings:

$$Corr(U_1, U_2) = \frac{\sum (r_{1,i} - \bar{r}_1)(r_{2,i} - \bar{r}_2)}{\sqrt{\sum (r_{1,i} - \bar{r}_1)^2} \sqrt{\sum (r_{2,i} - \bar{r}_2)^2}}$$

	L.N.	M.H.	J.H.	E.N.	D.U.	D.S.
E.N.	0.8704	-1.0000	NaN	NaN	NaN	-1.570092E-16

The single nearest student to E.N is L.N (Correlation : 0.8704)

b. Our answer is Python.

L.N is the nearest student to E.N.. Except for the common courses what L.N and E.N. complete, Python is the highest rating course in L.N's completed course.

c. cosine similarity between users

d.

Based on the cosine similarities, the nearest students to E.N is M.H. , J.H. and D.U. When we look up the common courses what they complete, Spatial is the only course which M.H. , J.H. and D.U. complete but E.N does not. That is the reason Spatial course should be recommended to E.N.

e.

兩者的差距在於，correlation有扣除平均值，而cosine則是直接拿原始的值來算相似度，而扣除平均值是為了把每個人評分的嚴謹程度算進去，像是對同一門課評價都為5分的人，其中可能包含：

1.給分嚴謹的人，他很少給出5分，這個5分的特殊性就較高。

2.給分彈性的人，他時常給出5分，這個5分的特殊性就較低。

透過扣掉平均，可以把這樣的狀況給算進去(第1種人扣掉平均有較高的數字)，而cosine則是不去扣掉平均，可以讓運算更加簡單，只看向量的關係性。

f.

i.

EN有修的課是SQL/DM.in.R/R.Prog/Regression

計算SQL/DM.in.R/R.Prog/Regression與其他課的相似度，並去預測EN對於其他課的喜好度，並且要有同時修兩堂課，都有評分的資料。

SQL + *Spatial *Python *Forecast

R.Prog + *Python *Forecast

Regression + *Python *Forecast

並且要找出有至少兩組評分的有兩組 (SQL, Spatial) (SQL, Python)並去計算其相似度。

ii.

推薦python由於EN有修的課，同樣修這些課的人，也有同樣修python並且評價不錯(3分/4分相比Forecast2/4分)

計算其相似度，透過相似度，推薦Spatial。

	row		column		cor
1		SQL	Spatial		0.8181
7		SQL	Python		-1.000

g.

透過item_based做預測最後得到的結果為下圖

	SQL	Spatial	PA.1	DM.in.R	Python	Forecast	R.Prog	Hadoop	Regression
[1,]	NA	4	2	3.26231	NA	NA	NA	NA	NA
[2,]	NA	NA	NA	3.00000	NA	3.000000	3	NA	3.000000
[3,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[4,]	NA	4	NA	NA	4	3.666667	NA	NA	NA
[5,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[6,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[7,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[8,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[9,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[10,]	NA	NA	NA	NA	2.000000	NA	NA	NA	NA
[11,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[12,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[13,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[14,]	NA	NA	NA	NA	NA	NA	NA	NA	NA
[15,]	NA	4	NA	NA	4	4.000000	NA	NA	3.296585

可以推薦 Spatial , Python , Forecast三門課

但是判斷R的算法可能與老師教得有些不同，由於Forecast無法計算相關性但卻出現在推薦並且根據上題 Spatial 及 Python 的相關性，推薦較大相關性的 Spatial 紿E.N.

Q3. Clustering

a.

- head(cluster1,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	War	Western
2	GoldenEye (1995)	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
17	From Dusk Till Dawn (1996)	0	1	0	0	0	0	1	1	0	0	0	0	1	0	0	0	1	0	0
21	Muppet Treasure Island (1996)	0	1	1	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0
24	Rumble in the Bronx (1995)	0	1	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
27	Bad Boys (1995)	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

- head(cluster2,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	war	western
14	Postino, Il (1994)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0
20	Angels and Insects (1995)	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
36	Mad Love (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0
51	Legends of the Fall (1994)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	1	1
55	Professional, The (1994)	0	0	0	0	0	0	1	0	0	1	0	0	0	0	1	0	1	0	0

- head(cluster3,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	war	western
6	Shanghai Triad (Yao a yao yao dao waipo qiao) (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
7	Twelve Monkeys (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0
9	Dead Man Walking (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
10	Richard III (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0
15	Mr. Holland's opus (1995)	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	1	0

- head(cluster4,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	war	western
3	Four Rooms (1995)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
5	Copycat (1995)	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	1	0	0
11	Seven (Se7en) (1995)	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0	0
12	Usual Suspects, The (1995)	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0	0
38	Net, The (1995)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0

- head(cluster5,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	war	western
1	Toy Story (1995)	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
4	Get Shorty (1995)	0	1	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0
8	Babe (1995)	0	0	0	0	0	1	1	0	0	1	0	0	0	0	0	0	0	0	0
13	Mighty Aphrodite (1995)	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
16	French Twist (Gazon maudit) (1995)	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0

- head(cluster6,n=5)

	Title	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	FilmNoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	war	western
32	Crumb (1994)	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
48	Hoop Dreams (1994)	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
71	Lion King, The (1994)	0	0	0	0	1	1	0	0	0	0	0	0	0	1	0	0	0	0	0
75	Brother Minister: The Assassination of Malcolm X (1994)	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
84	Robert A. Heinlein's The Puppet Masters (1994)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0

According to the chart above we can know that the first five movies in every cluster.

cluster1 : GoldenEye (1995)、From Dusk Till Dawn (1996)、Muppet Treasure Island (1996)、Rumble in the Bronx (1995)、Bad Boys (1995)

cluster2 : Postino, Il (1994)、Angels and Insects (1995)、Mad Love (1995)、Legends of the Fall (1994)、Professional, The (1994)

cluster3 : Shanghai Triad (Yao a yao yao dao waipo qiao) (1995)、Twelve Monkeys (1995)、Dead Man Walking (1995)、Richard III (1995)、Mr. Holland's Opus (1995)

cluster4 : Four Rooms (1995)、Copycat (1995)、Seven (Se7en) (1995)、Usual Suspects, The (1995)、Net, The (1995)

cluster5 : Toy Story (1995)、Get Shorty (1995)、Babe (1995)、Mighty Aphrodite (1995)、French Twist (Gazon maudit) (1995)

cluster6 : Crumb (1994)、Hoop Dreams (1994)、Lion King, The (1994)、Brother Minister: The Assassination of Malcolm X (1994)、Robert A. Heinlein's The Puppet Masters (1994)

b.

According to the chart above we can know that :

Cluster 1 can be classified as an action

Cluster 2 can be classified as an Romance :

Cluster 3 can be classified as an Drama :

Cluster 4 can be classified as an Thriller :

Cluster 5 can be classified as an Comedy and especially

Cluster 5 can be classified as an Conway and, especially, the average of Cluster 6 is too low to be classified into any category.

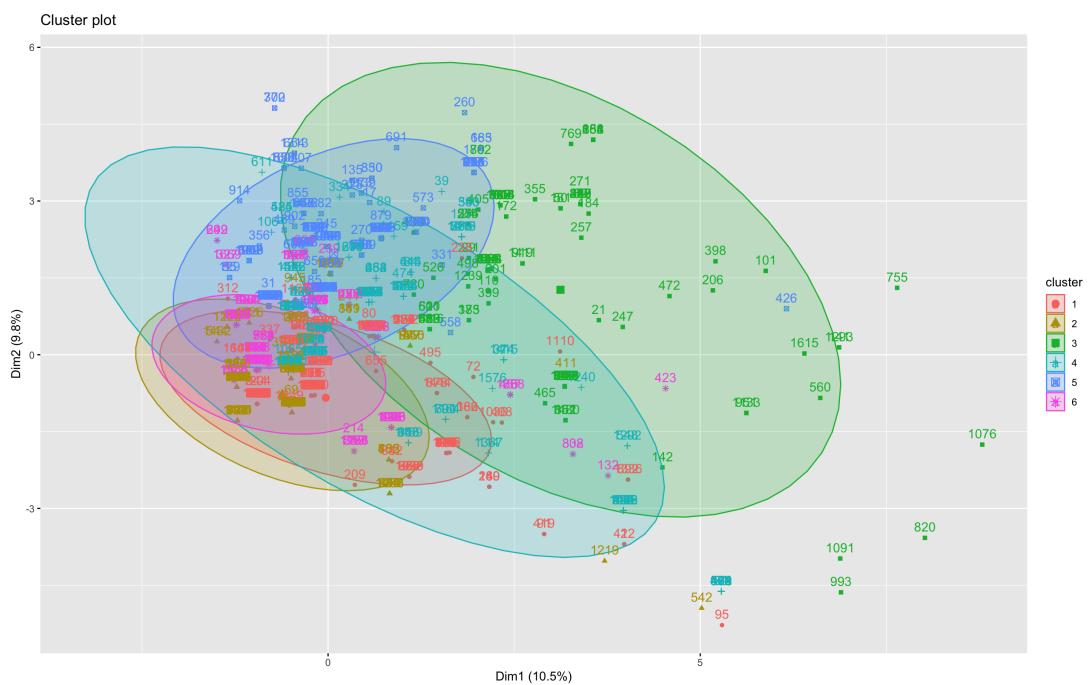
K-means clustering with 6 clusters of sizes 199, 100, 518, 118, 475, 254

Cluster means:	Unknown	Action	Adventure	Animation	Childrens	Comedy	Crime	Documentary	Drama	Fantasy	Filmnoir	Horror	Musical	Mystery	Romance	SciFi	Thriller	War	Western	
1	0.0000000000	1.0000000000	0.3610000000	0.0150757377	0.04252613	0.06521663	0.08547214	0.0000000000	0.05526738	0.020100503	0.0000000000	0.06523663	0.01005025	0.0317588	0.090442526	0.266311066	0.41708542	0.06030151	0.010050251	
2	0.0000000000	0.0700000000	0.0200000000	0.0000000000	0.1100000000	0.0400000000	0.0970000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.01158303	0.063706556	0.05212355	0.00772208	
3	0.0000000000	0.04826255	0.01737420	0.0000000000	0.03281853	0.0000000000	0.05984586	0.057391506	1.0000000000	0.007722008	0.003861004	0.00965251	0.01737452	0.01544402	0.07799661	0.11016944	0.2881356	0.06779961	0.11016949	0.00000000
4	0.0000000000	0.0000000000	0.0542373	0.0000000000	0.02542373	0.23728814	0.0000000000	0.01016942	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	0.0000000000	
5	0.0000000000	0.03789474	0.01684204	0.021052632	0.08421053	1.0000000000	0.02526362	0.02012053	0.16421053	0.041736842	0.0000000000	0.02315789	0.04210526	0.01473684	0.17473684	0.01894737	0.01894737	0.02315789	0.014736842	
6	0.007874016	0.0000000000	0.15354331	0.110236220	0.21259843	0.0000000000	0.06629913	0.181102362	0.0000000000	0.027559055	0.018509397	0.09824520	0.03149606	0.17473589	0.07086614	0.0000000000	0.0397300000	0.059118102		

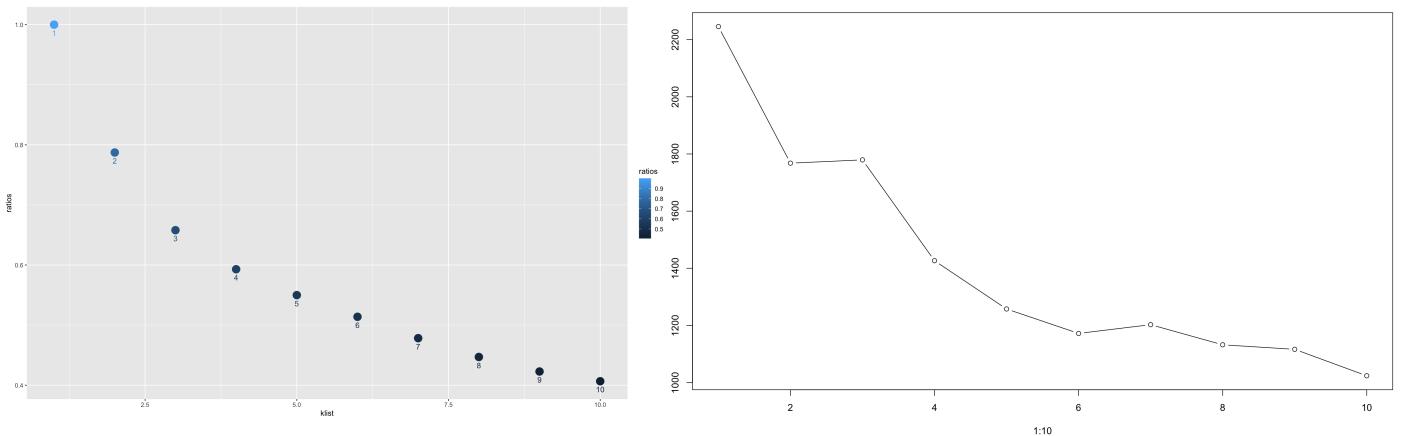
C.

When $k = 6$, the graph as shown below is the categorization results in part b.

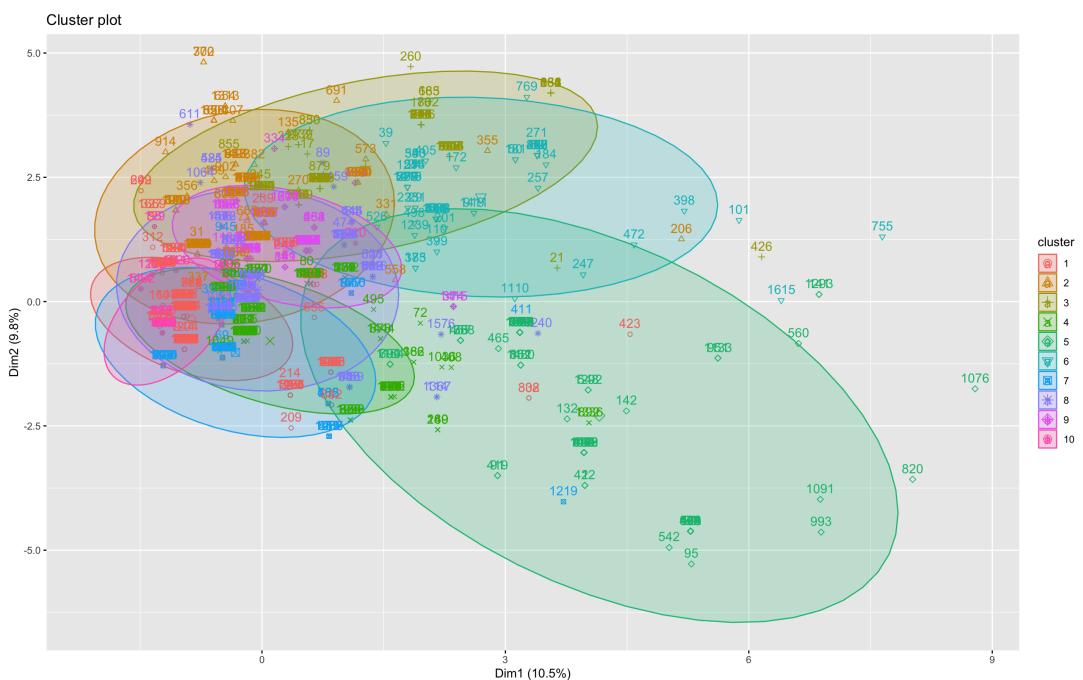
Obviously, there are still a lot of clusters overlapping. We will increase the number of k slowly to find the optimal solution.



When $k = 10$, the WSS value is lowest , and the ration about WSS(Within Cluster Sum of Squares)/TSS(Total Cluster Sum of Squares) is closest to 0. So we decide $k = 10$ to be the optimal solution.



It is much better than the result when setting $k = 6$.



Check number of movies in every cluster and make sure that there is no one cluster will has too few movies.

	cluster1	cluster2	cluster3	cluster4	cluster5	cluster6	cluster7	cluster8	cluster9	cluster10
number	529	151	84	304	71	68	130	165	78	84

When $k = 6$, the average of Cluster 6 is too low to be classified into any category. But when setting $k = 10$, We can observe obvious improvement. The number of movies which can not be categorized is decreased from 485 to 165, helping to make more movies be categorized precisely.