

**TRƯỜNG ĐẠI HỌC BÁCH KHOA
ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH**



BÁO CÁO BÀI TẬP NHÓM
Đề tài: Phân tích tập dữ liệu Las Vegas Strip

Môn học: Chuyên đề nghiên cứu khoa học dữ liệu ứng dụng
GVHD: PGS.TS Nguyễn Mạnh Tuân

Học viên:

Đinh Thanh Phong	2270243
Nguyễn Thị Hạnh Thảo	2270753
Nguyễn Công Trục	2270200
Hồ Duy Đạt Phúc	2051172

Mục lục

1. Đánh giá chất lượng bộ dữ liệu Lasvegas Tripadvisor:..	2
2. Mục tiêu phân tích:.....	3
3. Xử lý bộ dữ liệu sử dụng 3 giải thuật Rule Induction, Decision tree và k-NN.....	3
3.1) Giải thuật: Rule Induction.....	4
3.2) Giải thuật: k-NN.....	8
3.3) Giải thuật: Decision tree.....	11
4. Kết luận, đưa ra đánh giá về kết quả và đề xuất áp dụng mô hình cho nghiệp vụ khách sạn:.....	15

1.Đánh giá chất lượng bộ dữ liệu Lasvegas

Tripadvisor:

Bộ dữ liệu cung cấp các thông tin tổng kết từ những đánh giá trực tuyến của khách du lịch đã trải nghiệm thực tế và sử dụng dịch vụ tại các khách sạn nằm tại Lasvegas. Sự phát triển của ngành du lịch và các dịch vụ đặt phòng khiến cho sự cạnh tranh ngày càng gay gắt, ngành khách sạn phải đổi mới cách kinh doanh. Sự hài lòng và sự trung thành của khách hàng hết sức quan trọng đối với khách sạn. Do đó việc phân tích dữ liệu để hiểu được trải nghiệm của khách hàng với các dịch vụ được cung cấp trở thành một trong những xu hướng chính của lĩnh vực khách sạn. Từ đó đem đến góc nhìn tổng quát đa chiều cho người quản lý để đổi mới chất lượng các dịch vụ của khách sạn.

Nghiên cứu bộ dữ liệu có 504 dữ liệu đánh giá trực tuyến 21 khách sạn ở Las Vegas năm 2015 với các tiêu chí cụ thể bao gồm 20 thuộc tính: số trải nghiệm, mùa nghỉ dưỡng, điểm đánh giá, xếp hạng sao của khách sạn, cơ sở vật chất của khách sạn, các dịch vụ giải trí, thể thao của khách sạn, thông tin khách hàng như nhóm người thuộc quốc gia, khu vực, đối tượng khách hàng là doanh nhân, gia đình, bạn bè, độc thân hay những cặp đôi, số năm là thành viên, thời gian đánh giá,...Đối với các thuộc tính số phòng khách sạn, khách hàng tới từ khu vực, số năm là thành viên, ngày, tháng cung cấp đánh giá xuất hiện giá trị dữ liệu bị thiếu. Tuy nhiên vấn đề quan trọng của tập dữ liệu tập trung vào việc đánh giá chất lượng dịch vụ và khả năng đáp ứng dịch vụ của các khách

sạn thông qua các thuộc tính cụ thể có tính năng định lượng và phân loại từ các dữ liệu từ các nguồn dữ liệu chi tiết thu thập được từ người đánh giá như tiện ích bơi, phòng tập, tennis, spa, casino và truy cập internet miễn phí.

2. Mục tiêu phân tích:

Mục tiêu phân tích bộ dữ liệu nhằm nỗ lực khai thác thông tin để trả lời các câu hỏi nghiên cứu sau: Có thể dự đoán điểm số của một bài đánh giá khách sạn trực tuyến chỉ bằng cách sử dụng dữ liệu định lượng đầu vào không? Các tính năng ảnh hưởng đến hầu hết các điểm đánh giá trong khách sạn là gì? Làm thế nào để mỗi tính năng đó ảnh hưởng đến điểm số và kiến thức này có thể hữu ích cho các nhà quản lý khách sạn không? Bộ dữ liệu đã đóng góp các thông tin chính với mục tiêu:

- Tạo mô hình dự đoán điểm đánh giá dựa trên các đặc điểm định lượng của người dùng/người đánh giá và khách sạn, cũng như khoảng thời gian lưu trú cụ thể;
- Góp phần nghiên cứu phản hồi của khách hàng và đánh giá trực tuyến bằng cách cung cấp một cách tiếp cận mới về dữ liệu được sử dụng, các tính năng định lượng, trái ngược với các phân tích phổ biến nhất về chính văn bản đánh giá;
- Hiểu cách người dùng bị ảnh hưởng vốn có bởi các tính năng của khách sạn khi gửi điểm số bên cạnh nhận xét bằng văn bản trên các nền tảng trực tuyến, chẳng hạn như TripAdvisor.

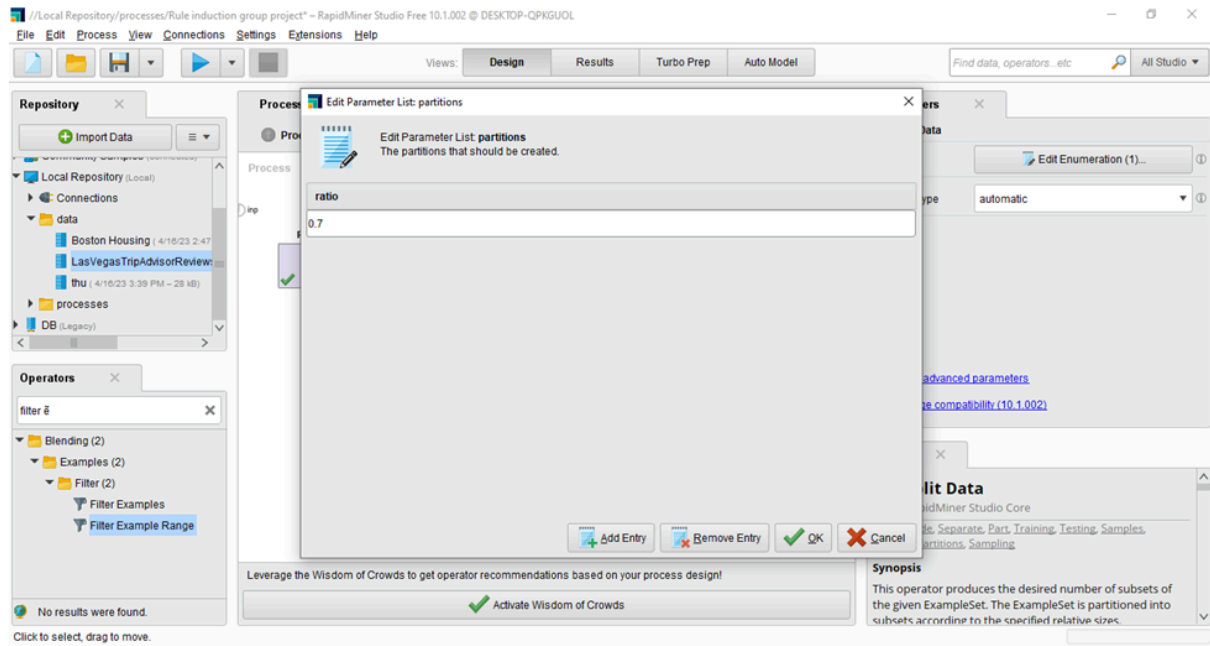
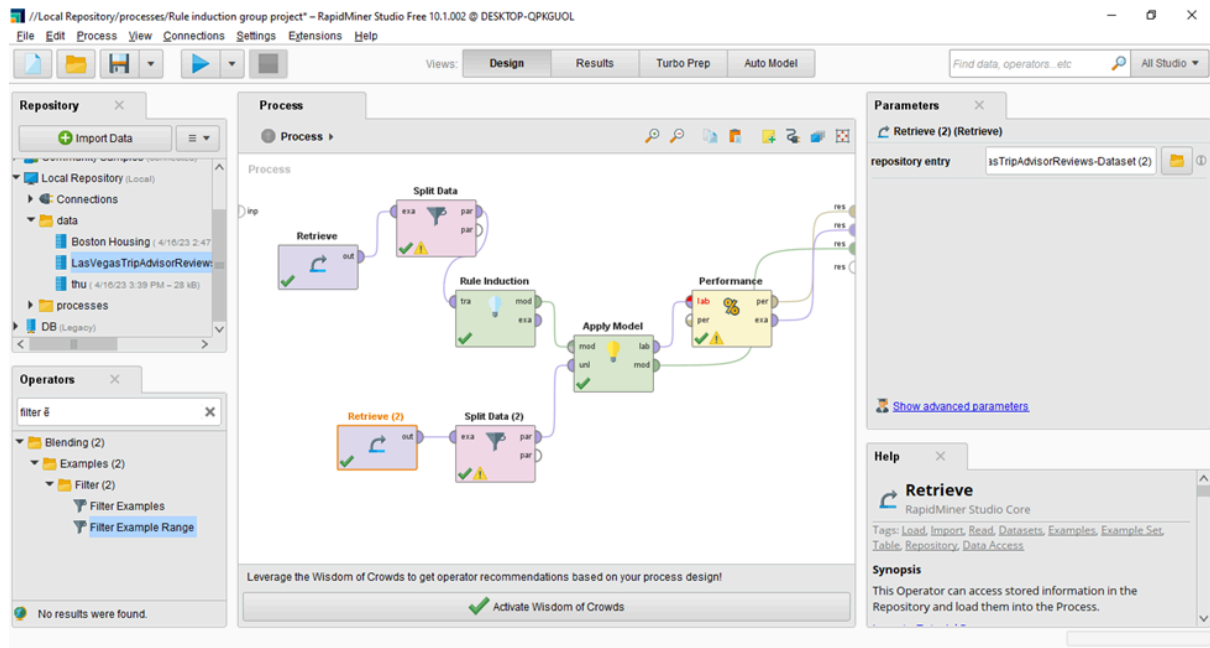
3. Xử lý bộ dữ liệu sử dụng 3 giải thuật Rule Induction, Decision tree và k-NN.

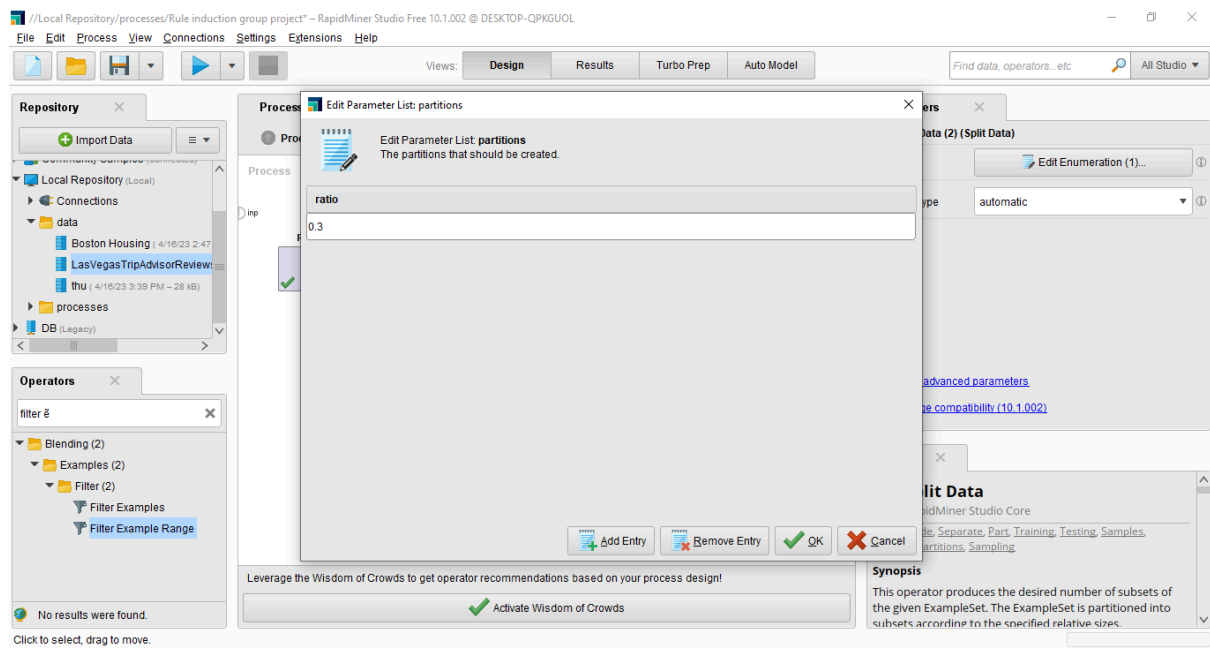
Nhóm dùng 3 nhiệm giải thuật Rule Induction, Decision tree và k-NN để thực hiện phân tích và so sánh kết quả (performance) của mô hình bằng cách sử dụng các công cụ sử dụng độ đo confusion matrix, precision, recall, accuracy.

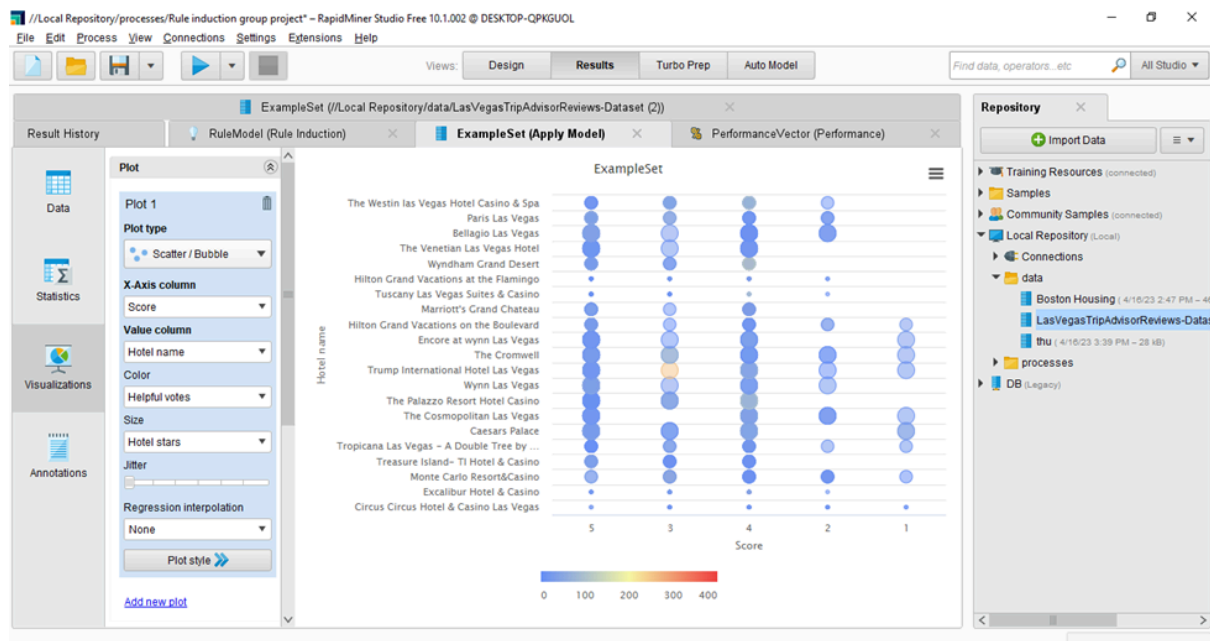
Để đạt được kết quả tốt nhất trong mỗi giải thuật, nhóm thực hiện lựa chọn các attributes khác nhau và chia tập train/test khác nhau ở mỗi giải thuật, nhằm đạt được kết quả tốt nhất cho bài toán.

- Thuộc tính “score” được chọn gán label.
- Sử dụng tất cả các thuộc tính khác làm attribute.
- Tỷ lệ phân chia Train/Test = 0.7/0.3









ExampleSet (/Local Repository/data/LasVegasTripAdvisorReviews-Dataset (2))

Table View

accuracy: 74.40%

	true 5	true 3	true 4	true 2	true 1	class precision
pred. 5	187	8	20	10	4	81.66%
pred. 3	1	38	3	0	0	90.48%
pred. 4	35	21	135	7	0	68.18%
pred. 2	4	5	6	13	5	39.39%
pred. 1	0	0	0	0	2	100.00%
class recall	82.38%	52.78%	82.32%	43.33%	18.18%	

Kết quả giải thuật Rule Induction:

- Recall, Precision: như bảng đính kèm ở trên.
- Accuracy: 74,40%

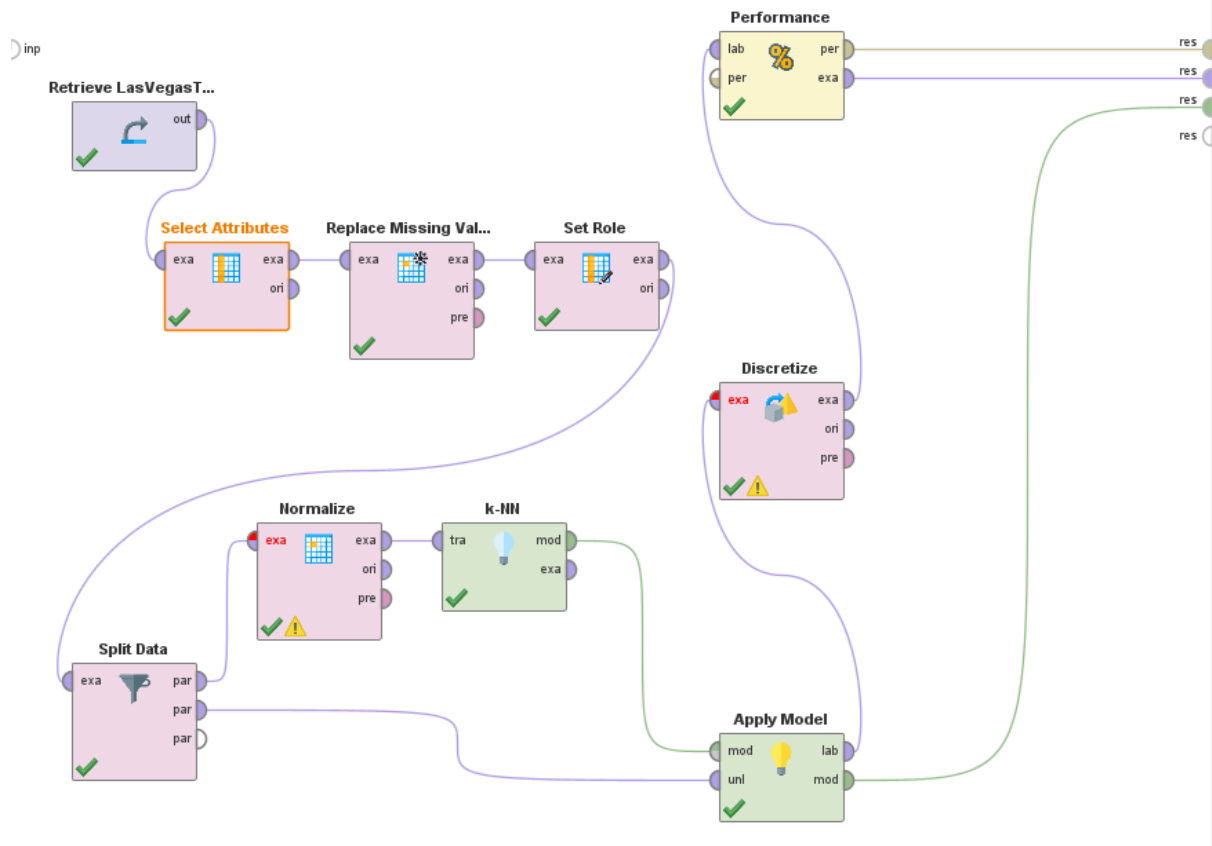
3.2) Giải thuật: k-NN.

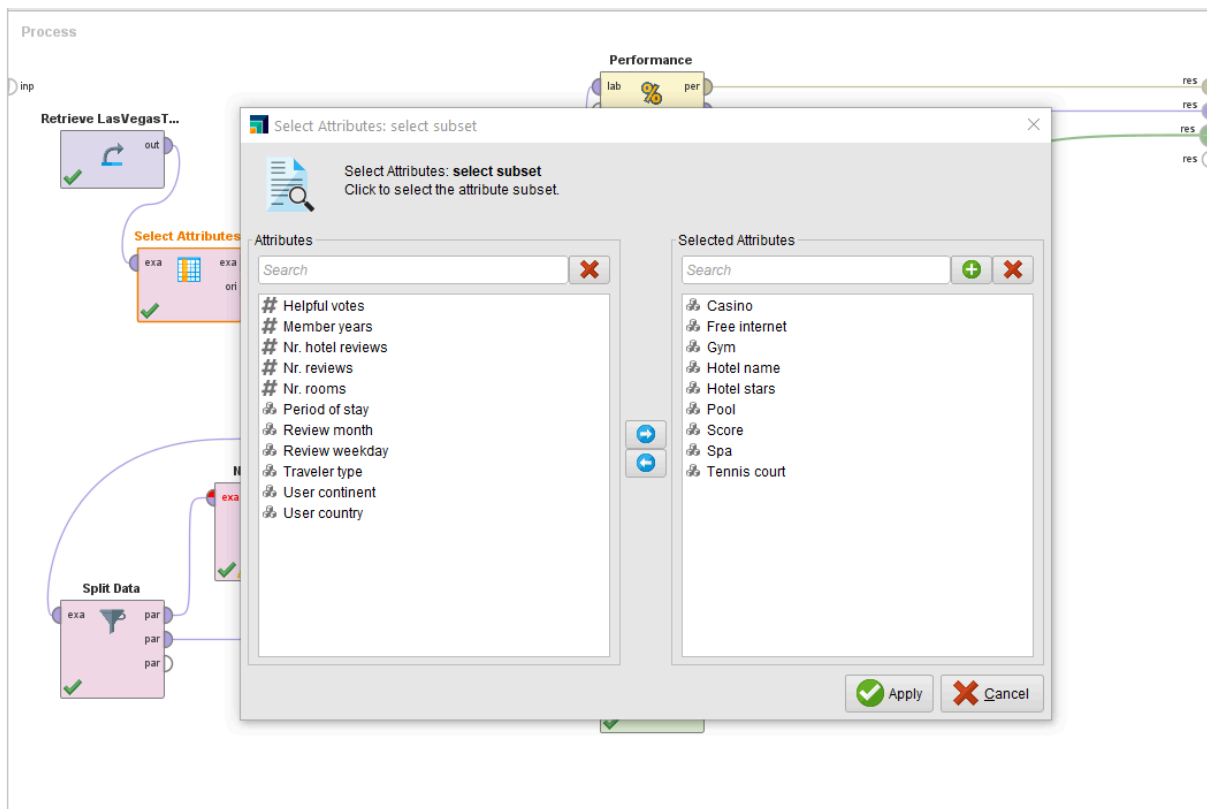
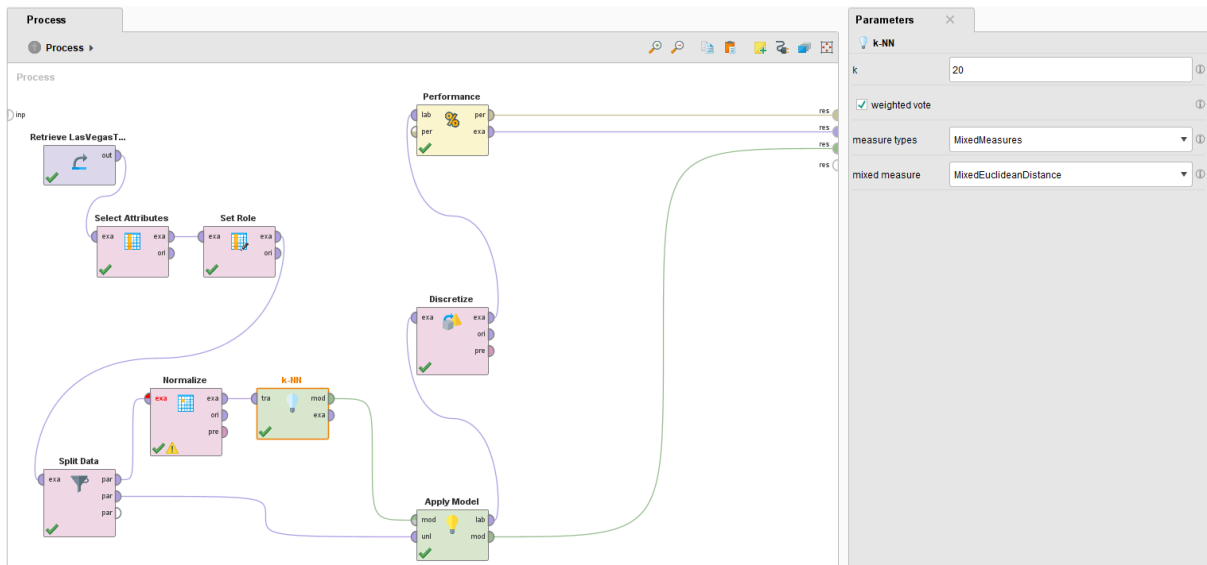
- Thuộc tính “score” được chọn gắn label.

- Các thuộc tính được chọn làm attribute: Casino, Free Internet, Gym, Hotel name, Hotel star, Pool, Spa, Tennis court.
- Chọn $k = 20$
- Tỷ lệ phân chia Train/Test = 0.8/0.2

Row No.	Score	User country	Nr. reviews	Nr. hotel revi...	Helpful votes	Period of stay	Traveler type	Pool	Gym	Tennis court	Spa	Casino	Free internet	Hotel name	Hotel stars	Nr. rooms	User conti
1	5	USA	11	4	13	Dec-Feb	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
2	3	USA	119	21	75	Dec-Feb	Business	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
3	5	USA	36	9	25	Mar-May	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
4	4	UK	14	7	14	Mar-May	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Europe
5	4	Canada	5	5	2	Mar-May	Solo	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
6	3	Canada	31	8	27	Mar-May	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
7	4	UK	45	12	46	Mar-May	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Europe
8	4	USA	2	1	4	Mar-May	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
9	4	India	24	3	8	Mar-May	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Asia
10	3	Canada	12	7	11	Mar-May	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
11	2	USA	102	24	58	Jun-Aug	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
12	3	Australia	20	9	24	Jun-Aug	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Oceania
13	2	USA	7	6	9	Jun-Aug	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
14	3	USA	22	5	13	Jun-Aug	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
15	3	UK	3	3	0	Jun-Aug	Friends	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Europe
16	4	New Zealand	146	17	33	Jun-Aug	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Oceania
17	1	Canada	8	8	9	Sep-Nov	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
18	4	USA	9	3	1	Sep-Nov	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
19	3	Canada	41	9	19	Sep-Nov	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
20	2	USA	8	7	26	Sep-Nov	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer
21	4	UK	10	5	2	Sep-Nov	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Europe
22	1	New Zealand	4	3	3	Sep-Nov	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Oceania
23	4	UK	18	7	19	Dec-Feb	Families	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	Europe
24	2	USA	4	4	3	Dec-Feb	Couples	NO	YES	NO	NO	YES	YES	Circus Circus...	3	3773	North Amer

Process





The screenshot shows the Orange3 interface. On the left, a 'Split Data' widget is connected to an 'ExampleSet (Discretize)' widget. The 'Edit Parameter List: partitions' dialog is open, showing a 'ratio' parameter with values 0.8 and 0.2. Below the dialog, the 'PerformanceVector (Performance)' widget is displayed in 'Table View'.

Edit Parameter List: partitions

ratio

0.8

0.2

Buttons: Add Entry, Remove Entry, OK, Cancel

PerformanceVector (Performance)

Table View

accuracy: 55.00%

	true 5	true 3	true 4	true 2	true 1	class precision
pred. 5	41	7	17	3	2	58.57%
pred. 3	0	2	4	0	0	33.33%
pred. 4	4	5	12	3	0	50.00%
pred. 2	0	0	0	0	0	0.00%
pred. 1	0	0	0	0	0	0.00%
class recall	91.11%	14.29%	36.36%	0.00%	0.00%	

Kết quả giải thuật k-NN:

- Recall, Precision: như bảng đính kèm ở trên.
- Accuracy: 55,50%

3.3) Giải thuật: Decision tree.

- Thuộc tính “score” được chọn gắn label.
- Các thuộc tính được chọn làm attribute: Casino, Helpful votes, Nr. reviews, Nr. rooms, Period of stay, Review months, Review weekday, Travel type, User continent, User country.

- Tỷ lệ phân chia Train/Test = 0.7/0.3

Local Repository\data/decisiontree - RapidMiner Studio Free 10.1.001 @ MSI

File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep Auto Model

Find data, operators, etc. All Studio

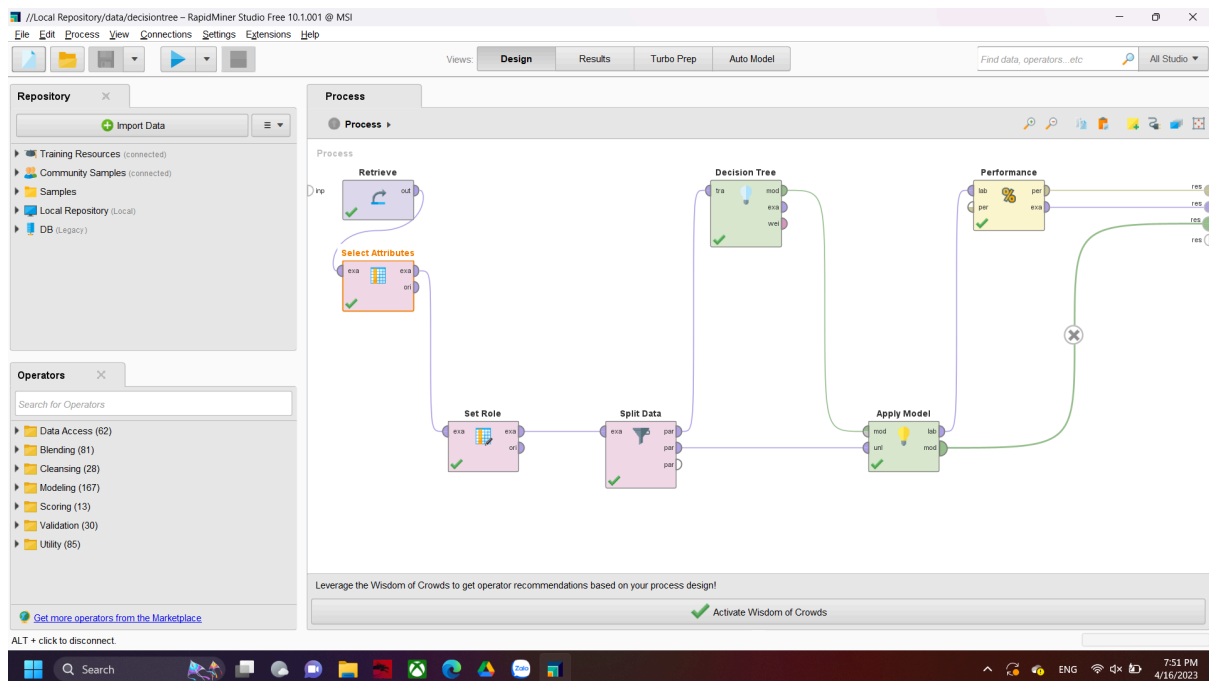
Result History: ExampleSet (Retrieve) Tree (Decision Tree) ExampleSet (Apply Model) PerformanceVector (Performance)

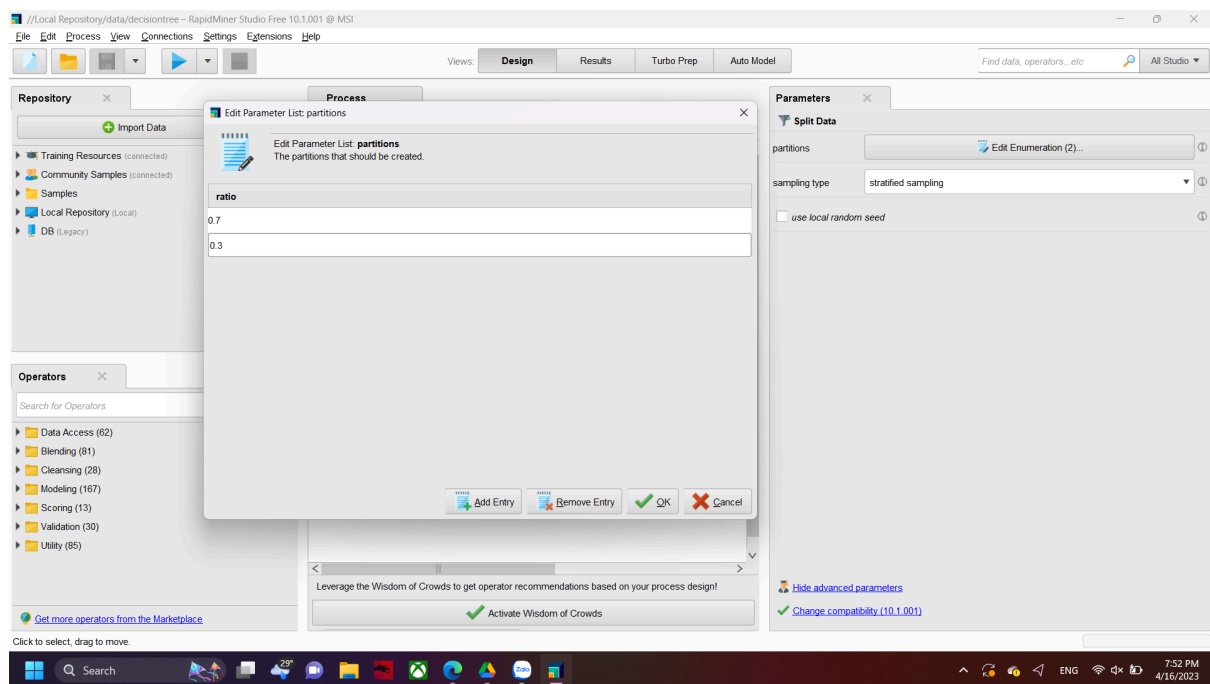
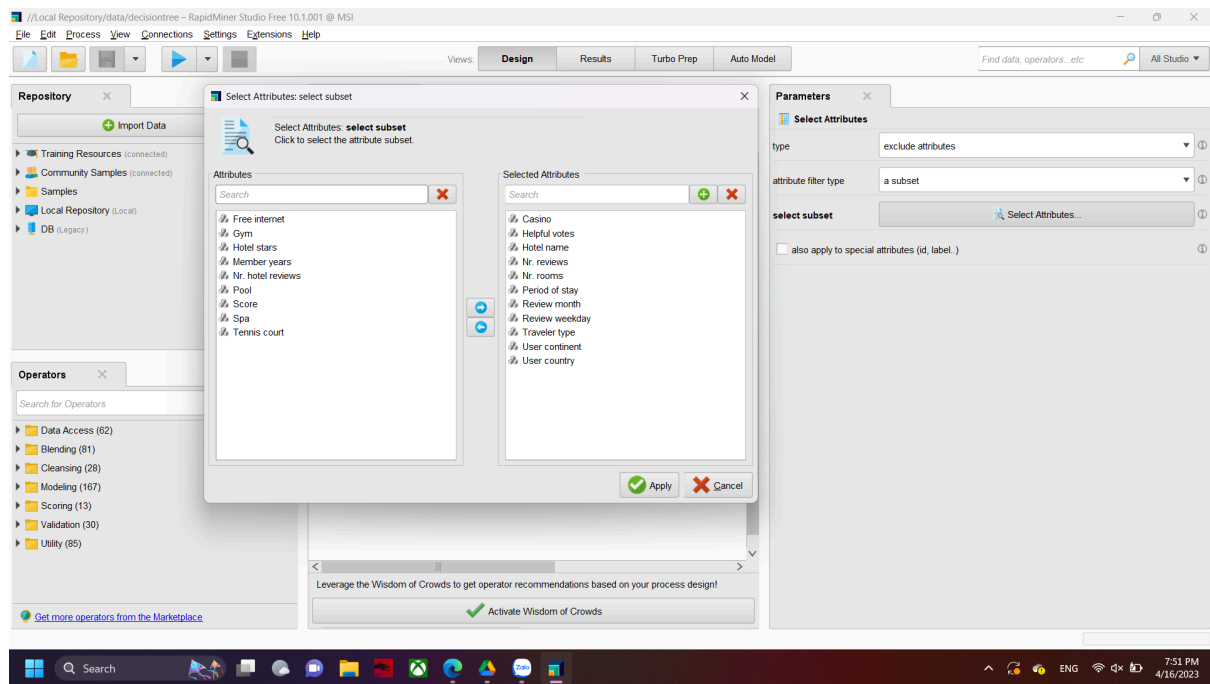
Open in: Turbo Prep Auto Model

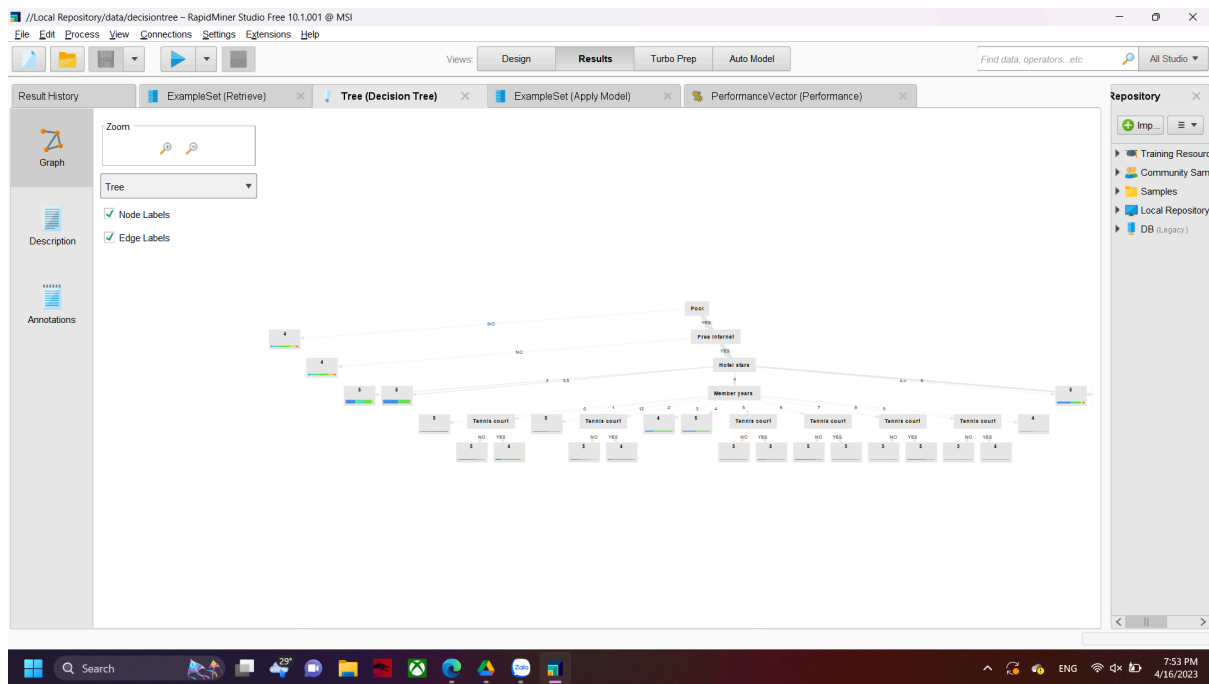
Filter (504 / 504 examples): all

Row No.	User country	Nr. reviews	Nr. hotel rev...	Helpful votes	Score	Period of stay	Traveler type	Pool	Gym	Tennis court	Spa	Casino
1	USA	11	4	13	5	Dec-Feb	Friends	NO	YES	NO	NO	YES
2	USA	119	21	75	3	Dec-Feb	Business	NO	YES	NO	NO	YES
3	USA	36	9	25	5	Mar-May	Families	NO	YES	NO	NO	YES
4	UK	14	7	14	4	Mar-May	Friends	NO	YES	NO	NO	YES
5	Canada	5	5	2	4	Mar-May	Solo	NO	YES	NO	NO	YES
6	Canada	31	8	27	3	Mar-May	Couples	NO	YES	NO	NO	YES
7	UK	45	12	46	4	Mar-May	Couples	NO	YES	NO	NO	YES
8	USA	2	1	4	4	Mar-May	Families	NO	YES	NO	NO	YES
9	India	24	3	8	4	Mar-May	Friends	NO	YES	NO	NO	YES
10	Canada	12	7	11	3	Mar-May	Families	NO	YES	NO	NO	YES
11	USA	102	24	58	2	Jun-Aug	Families	NO	YES	NO	NO	YES
12	Australia	20	9	24	3	Jun-Aug	Friends	NO	YES	NO	NO	YES
13	USA	7	6	9	2	Jun-Aug	Friends	NO	YES	NO	NO	YES
14	USA	22	5	13	3	Jun-Aug	Friends	NO	YES	NO	NO	YES
15	UK	3	3	0	3	Jun-Aug	Friends	NO	YES	NO	NO	YES
16	New Zealand	146	17	33	4	Jun-Aug	Families	NO	YES	NO	NO	YES
17	Canada	8	8	9	1	Sep-Nov	Families	NO	YES	NO	NO	YES

ExampleSet (504 examples, 0 special attributes, 20 regular attributes)







The screenshot shows the RapidMiner Studio interface. The top menu bar includes File, Edit, Process, View, Connections, Settings, Extensions, and Help. Below the menu is a toolbar with icons for file operations and a search bar. The main workspace contains several tabs: Result History, ExampleSet (Retrieve), Tree (Decision Tree), ExampleSet (Apply Model) (selected), and PerformanceVector (Performance). On the left sidebar, there are icons for Data, Statistics, Visualizations, and Annotations. The central area displays a table titled "ExampleSet (Apply Model)" with 17 rows of data. The table has columns for Row No., Score, prediction(Score), confidence(5), confidence(3), confidence(4), confidence(2), confidence(1), Nr. hotel rev..., Pool, Gym, Tennis court, and Spa. The data is filtered to show 151 examples.

Row No.	Score	prediction(Score)	confidence(5)	confidence(3)	confidence(4)	confidence(2)	confidence(1)	Nr. hotel rev...	Pool	Gym	Tennis court	Spa
1	3	4	0.067	0.267	0.400	0.200	0.067	21	NO	YES	NO	NO
2	5	4	0.067	0.267	0.400	0.200	0.067	9	NO	YES	NO	NO
3	4	4	0.067	0.267	0.400	0.200	0.067	5	NO	YES	NO	NO
4	2	4	0.067	0.267	0.400	0.200	0.067	24	NO	YES	NO	NO
5	3	4	0.067	0.267	0.400	0.200	0.067	9	NO	YES	NO	NO
6	3	4	0.067	0.267	0.400	0.200	0.067	5	NO	YES	NO	NO
7	4	4	0.067	0.267	0.400	0.200	0.067	3	NO	YES	NO	NO
8	1	4	0.067	0.267	0.400	0.200	0.067	3	NO	YES	NO	NO
9	4	4	0.067	0.267	0.400	0.200	0.067	7	NO	YES	NO	NO
10	4	5	0.347	0.327	0.265	0.061	0	42	YES	YES	NO	YES
11	3	5	0.347	0.327	0.265	0.061	0	3	YES	YES	NO	YES
12	3	5	0.347	0.327	0.265	0.061	0	8	YES	YES	NO	YES
13	4	5	0.347	0.327	0.265	0.061	0	27	YES	YES	NO	YES
14	4	5	0.347	0.327	0.265	0.061	0	263	YES	YES	NO	YES
15	4	5	0.347	0.327	0.265	0.061	0	6	YES	YES	NO	YES
16	4	5	0.347	0.327	0.265	0.061	0	8	YES	YES	NO	YES
17	4	4	0.062	0.250	0.438	0.188	0.062	8	YES	YES	NO	YES

accuracy: 52.32%

	true 5	true 3	true 4	true 2	true 1	class precision
pred. 5	65	14	35	6	2	53.28%
pred. 3	0	1	1	0	0	50.00%
pred. 4	3	7	13	3	1	48.15%
pred. 2	0	0	0	0	0	0.00%
pred. 1	0	0	0	0	0	0.00%
class recall	95.59%	4.55%	26.53%	0.00%	0.00%	

Kết quả giải thuật Decision tree:

- Recall, Precision: như bảng đính kèm ở trên.
- Accuracy: 55,50%

4. Kết luận, đưa ra đánh giá về kết quả và đề xuất áp dụng mô hình cho nghiệp vụ khách sạn:

Thông qua việc sử dụng 3 giải thuật: Rule induction, Decision tree và k-NN nhận thấy giải thuật Rule induction cho độ chính xác (*accuracy*) cao nhất 74,4% so với giải thuật decision tree có độ chính xác 52,32% và giải thuật k-NN có độ chính xác 55%.

Từ model xây dựng được qua các giải thuật ở trên, người dùng có thể dự đoán được “score” cho các hotel mà khách đánh giá. Người dùng có một cái nhìn tổng quan về mức quan trọng và các yếu tố ảnh hưởng đến

“score”. Từ đó, người dùng có thể áp dụng các nghiệp vụ chuyên ngành để cải thiện dịch vụ và nâng cao chất lượng dịch vụ hơn.

