Module 8

8.4 Unsupervised learning: clustering and association rules in KNIME

By María Martínez Rojas

Associate Professor CD, University of Granada

By José Manuel Soto Hidalgo

Associate Professor, CEAR, University de Granada

This capsule focuses on how to implement the different unsupervised learning algorithms introduced in Module 6 (*Unsupervised learning: clustering and association rules*) in *KNIME*. We will create data flows representative of the data science life cycle to solve clustering problems and create association rules. As in the previous capsule, data sets already used in Modules 2, 3, and 6 (*Bioinformatic analysis of an omics problem, Data science and machine learning,* and *Unsupervised learning: clustering and association rules*) will be used as examples here. Specifically, we will analyze the gene expression data discussed in Module 2 and used as example data in Module 6.

1.1. Clustering

In this section we will create data flows to solve a clustering problem which will specifically focus on identifying groups in the data without using any a priori information known about the categories, types, classes, or groups in the data. The flow represented in figure 1 shows an implementation of the two types of clustering discussed in Module 6: hierarchical and *k*-means clustering. Hierarchical clustering can be implemented in two different ways, either by including an external distance model and with its own node, or by including the distance matrix. Finally, the results are visualized through a heatmap and different visualization options (continuous or discrete) are shown.

1









Figure 1. Data flow with hierarchical clustering and k-means analysis.

First, we read the data (expression matrix) with the "File Reader" node and, as an example used to facilitate our understanding of the dendrograms, we will select only a portion of the expression matrix (50 genes and 20 samples) with the "Column Filter" and "Row Filter" nodes. Once the 50 genes and 20 samples are selected, we demonstrate several ways to perform hierarchical clustering with *KNIME*. A distance model can be created ("Numeric Distances" node) which is used by the "Hierarchical Clustering" node ("DistMatrix") together with the data to perform the hierarchical clustering. Figure 2 shows the different configuration options of the "Numeric Distances" node in which different distances can be used (e.g., Euclidean, Manhattan, or Maximo, etc.) the variables (samples) to be considered in the clustering, and treatment of missing values.







	Dialog - 4:31 - Nume	ric Distances (MODELO DE)
	Distance Configuration	Flow Variables Memory Policy
Exclude	• Manual Selection	Wildcard/Regex Selection Include T Filter D TCGA-EE-A2MC-06 D TCGA-EE-A2MF-06 D TCGA-EE-A2MF-06 D TCGA-EE-A2MF-06 D TCGA-EE-A2MF-06 D TCGA-EE-A2MF-06 D TCGA-EE-A2MJ-06
Enforce exclusion istance Selection Standard Distance	(Euclidean/ Manhattan)	Enforce inclusion
Configuration Euclidean Manhattan		
Maximum Custom 'p' 2.0		
Normalize distance	(Requires normalized input vectors)	
Missing Values	(fails if a missing value cell occurs during computa (Assume a missing value has the value of the respe	iion.) ctive counterpart – this will add 0 to the sum of pair wise absolute differences)

Figure 2. Options and parameters of the "Numeric Distances" node.

The "Hierarchical Clustering" node ("DistMatrix") generates a data model with the clusters that, together with the "Hierarchical Cluster Assigner" node and the data, assigns each entry to a cluster number. Finally, we can draw the heatmap with the "Heatmap" node and interact with it, for example, by visualizing it with continuous (figure 3) or discrete (figure 4) representation.









Figure 3. Heatmap visualization of results with continuous representation.



Figure 4. Heatmap results displayed as a discrete representation.



UNIVERSIDAD DE GRANADA



abiertaugi

Hierarchical clustering can also be performed without using a distance model by using the "Hierarchical Clustering" node. This node presents options for the type of distance to use, number of output clusters, and type of union, as well as the variables to consider (figure 5).

	Options Flow Variables	Memory Policy
	Options Flow variables	Memory Foncy
	Number output cluster:	3
	Distance function: Euc	lidean 😴
	Linkara taran SINCI	
	Linkage type: SINGL	
	Cache distan	605
		les
Exclude	r	Include
T Filter		T Filter
1 mer		
	>	D TCGA-D9-A426-06
		D TCGA-EE-A3AF-06
	>>>	D TCGA-ER-A19F-06
		D TCGA-EE-A2MF-06
	<	D TCGA-EE-A2MJ-06
		D TCCA-BF-AAP4-01 D TCCA-D3-A8CM-06
		$\Box \Box \Box \Box C G \Delta - G N - \Delta 26 \Delta - 06$

Figure 5. The "Hierarchical Clustering" node options and parameters.

We can also visualize and interact with the dendrogram by employing this node. To do so, right click on the node once it has been executed and select the "View: Dendrogram"/"distance view" option to show the result (figure 6).









Figure 6. Visualization of the results as a dendrogram.

Finally, the "*k*-means clustering" option is also shown in the flow shown in figure 1; the "*k*-Means" node allows us to perform such clustering in a very simple way. As shown in figure 7, it allows us to set the number of clusters to obtain and randomly initiate a maximum number of iterations, etc.







Clusters Number of clusters: First k rows Random initialization Max. number of iterations: Solumn Selection Exclude Filter	ndom seed 0 New
Clusters Number of clusters: Centroid initialization: First k rows Random initialization Use static ra Number of Iterations Max. number of iterations: Scolumn Selection Exclude Filter	ndom seed 0 New
Number of clusters: 3 2 . Centroid initialization: First k rows Random initialization Vuse static ra Number of Iterations Max. number of iterations: 99 Column Selection Exclude Filter	ndom seed 0 New
Centroid initialization: First k rows Random initialization Vuse static ra Number of Iterations Max. number of iterations: Column Selection Exclude Filter	New
 First k rows Random initialization	New
● Random initialization ✓ Use static ra Number of Iterations ● Max. number of iterations: ● Column Selection ● Exclude ● T Filter ●	New
 Random initialization Use static ra Number of Iterations Max. number of iterations: 99 Column Selection Exclude Filter 	New
Number of Iterations Max. number of iterations: 99 Column Selection Exclude Filter	 Include Filter D TCGA-D9-A4Z6-06 D TCGA-EE-A2MQ-06 D TCGA-EE-A3AF-06 D TCGA-ER-A19F-06 D TCGA-ER-A19F-06
Max. number of iterations: 99 Column Selection Exclude Filter	 Include Filter D TCGA-D9-A4Z6-06 D TCGA-EE-A2MQ-06 D TCGA-EE-A3AF-06 D TCGA-ER-A19F-06 D TCGA-ER-A19F-06
Column Selection Exclude Filter	 Include <i>Filter</i> D TCGA-D9-A4Z6-06 D TCGA-EE-A2MQ-06 D TCGA-EE-A3AF-06 D TCGA-ER-A19F-06 D TCGA-ER-A19F-06
Column Selection Exclude Filter	 Include <i>Filter</i> D TCGA-D9-A4Z6-06 D TCGA-EE-A2MQ-06 D TCGA-EE-A3AF-06 D TCGA-ER-A19F-06 D TCGA-ER-A19F-06
	 CGA-EE-A2MF-06 TCGA-EE-A2MJ-06 TCGA-BF-AAP4-01 TCGA-D3-A8GM-06 TCGA-GN-A26A-06 TCGA-EB-A3XE-01
	Always include all columns
lilite Mapping	
Z Enable Hilite Mapping	

Figure 7. The options and parameters for the "k-Means" node.

1.2. Association rules

In this section we will create data flows to solve an association rules problem (figure 8) using the same data set as used in Module 6 (Capsule 3, *Association rules*), which comprises six variables:

7

- MUTATIONSUBTYPES
- UV-signature
- RNASEQ-CLUSTER_CONSENHIER
- MethTypes.201408
- MIRCluster

UNIVERSIDAD

DE GRANADA

LYMPHOCYTE.SCORE







Figure 8. Data flow for association rules.

First, we use the "Excel Reader" node to read the data set file containing all the variables in the set. The "Column Filter" node is used to select 6 of the variables through the configuration options of this node, as shown in figure 9.

Column Filter	low Variable	es Memory Policy
Manual Selection Will Will Will Filter ALL_PRIMARY_VS_METASTATIC ALL_PRIMARY_VS_METASTATIC REGIONAL_VS_PRIMARY ProteinCluster OncoSignCluster OncoSignCluster CDKN2A(cg13601799)_meth KIT(cg10087973)_meth BRAF_cna NRAS_cna	dcard/Rege	x Selection Type Selection Include Filter S MUTATIONSUBTYPES UV-signature RNASEQ-CLUSTER_CONSENHIER MethTypes.201408 MIRCluster LYMPHOCYTE.SCORE
Enforce exclusion	-	Enforce inclusion

Figure 9. The "Column Filter" node configuration menu.

These two steps conclude the data preprocessing required to perform the transactions and calculate the association rules. As explained in Module 6 (*Unsupervised learning: clustering and association rules*), the first step in obtaining association rules is identifying what defines the data items and transactions. To do this we will use the "Column





Combiner" and "Cell Splitter" nodes to join the values of the variables into a single column and group them into an array. Figure 10 shows the output of the "Cell Splitter" node, showing the newly created "Transaction" column containing all the values for each of the variables in an array.

		Output Table - 3:29 - Cell Splitter (TF	RANSFORMA)
File Hilite	Navigation	View	
	(Table "default" - Rows: 331 Spec - Columns: 8	Properties Flow Variables
Row ID	S LYMP	S Transaccion] Transaccion_SplitResultList
Row0	2	"BRAF_Hotspot_Mutants","UV signature","keratin","norma.	. ["BRAF_Hotspot_Mutants","UV signature","keratin",]
Row1	4	"RAS_Hotspot_Mutants","UV signature","keratin","CpG isl.	["RAS_Hotspot_Mutants","UV signature","keratin",]
Row2	5	"BRAF_Hotspot_Mutants","UV signature","keratin","norma.	. ["BRAF_Hotspot_Mutants","UV signature","keratin",]
Row3	2	"RAS_Hotspot_Mutants","UV signature","keratin","hypo-m	. ["RAS_Hotspot_Mutants","UV signature","keratin",]
Row4	6	"Triple_WT", "not UV", "immune", "CpG island-methylated".	. ["Triple_WT","not UV","immune",]
Row5	4	"BRAF_Hotspot_Mutants","UV signature","keratin","hypo	["BRAF_Hotspot_Mutants","UV signature","keratin",]
Row6	0	"BRAF_Hotspot_Mutants","UV signature","keratin","norma.	. ["BRAF_Hotspot_Mutants","UV signature","keratin",]
Row7	0	"BRAF_Hotspot_Mutants","UV signature","MITF-low","hyp.	["BRAF_Hotspot_Mutants","UV signature","MITF-low",]
Row8	6	"BRAF_Hotspot_Mutants","UV signature","MITF-low","hyp.	["BRAF_Hotspot_Mutants","UV signature","MITF-low",]
Row9	5	"RAS_Hotspot_Mutants","UV signature","keratin","hypo-m	. ["RAS_Hotspot_Mutants","UV signature","keratin",]
Row10	5	"-","-","keratin","hypo-methylated","MIR.type.2","5"	["-","-","keratin",]
Row11	5	"BRAF_Hotspot_Mutants","UV signature","immune","CpG i	. ["BRAF_Hotspot_Mutants","UV signature","immune",]
Row12	4	"-","-","keratin","CpG island-methylated","MIR.type.3","4	["-","-","keratin",]
Row13	3	"-","-","immune","CpG island-methylated","MIR.type.4","	" ["-","-","immune",]
Row14	2	"RAS_Hotspot_Mutants","not UV","keratin","hyper-methyl.	. ["RAS_Hotspot_Mutants","not UV","keratin",]
Row15	6	"Triple_WT","not UV","immune","CpG island-methylated".	. ["Triple_WT","not UV","immune",]
Row16	5	"BRAF_Hotspot_Mutants","UV signature","MITF-low","hyp.	["BRAF_Hotspot_Mutants","UV signature","MITF-low",]
Row17	2	"BRAF_Hotspot_Mutants","UV signature","MITF-low","hyp.	["BRAF_Hotspot_Mutants","UV signature","MITF-low",]
Row18	6	"BRAF_Hotspot_Mutants","UV signature","MITF-low","hyp.	["BRAF_Hotspot_Mutants","UV signature","MITF-low",]
Row19	6	"BRAF_Hotspot_Mutants","UV signature","immune","hyper	. ["BRAF_Hotspot_Mutants","UV signature","immune",]
Row20	5	"RAS_Hotspot_Mutants","UV signature","MITF-low","hypo.	["RAS_Hotspot_Mutants","UV signature","MITF-low",]
Row21	5	"BRAF_Hotspot_Mutants","not UV","immune","hypo-meth.	["BRAF_Hotspot_Mutants","not UV","immune",]
Row22	2	"BRAF_Hotspot_Mutants","not UV","immune","CpG island.	["BRAF_Hotspot_Mutants","not UV","immune",]
Row23	5	"RAS_Hotspot_Mutants","UV signature","immune","normal	. ["RAS_Hotspot_Mutants","UV signature","immune",]
Row24	6	"Triple_WT","not UV","keratin","normal-like","MIR.type.2".	. ["Triple_WT","not UV","keratin",]
Row25	0	"BRAF_Hotspot_Mutants","UV signature","immune","hyper	. ["BRAF_Hotspot_Mutants","UV signature","immune",]
Row26	5	"RAS_Hotspot_Mutants","UV signature","immune","normal	. ["RAS_Hotspot_Mutants","UV signature","immune",]
Row27	5	"Triple WT", "UV signature", "immune", "normal-like", "MIR	. ["Triple WT","UV signature","immune",]

Figure 10. Output of the "Cell Splitter" node showing the transactions as an array.

These transactions can be used to determine a set of frequent items and to extract association rules. In the following, we illustrate how to represent these two examples in *KNIME*.

1. Determine a set of frequent item set

We can search for frequent items in a list of item sets with the "Item Set Finder" node which provides different algorithms for this task, including "A priori", "FP-growth", "RElim", "Sam", "JIM", "DICE", "TANIMOTO", as shown in figure 11. In addition, we can also determine the objective ("Frequent", "Closed", or "Maximal") and assign the Minimum set size and Minimum support levels for the items. In this example, we will explore the "A priori" algorithm (as detailed in Module 6, *Unsupervised learning: clustering and association rules*) which has a support value of 0.015. Figure 12 shows the item sets obtained when applying the "Item Set Finder" node with these settings, sorted by the relative percentage of support.





D	ialog - 0:32 - Item Set	Finder (Borgelt) (IT	EMSET)
Options	Advanced Settings	Flow Variables	Memory Policy
lte	m column: [] Trans	accion_SplitResult	List ᅌ
Algortihm: • Apriori I	Pgrowth 🔵 RElim 🤇	SaM JIM	
	Target Type:	Closed 🔵 Maxim	al
Item set settings			
	Minimum set size:		1 0
	Minimum support:	0.01	5 0
	 Absolute num 	iber 🔵 Percentag	e
Threshold: (opt	ional) 10.0 0		Sort item set

Figure 11. Options and parameters for the "Item Set Finder" node.

D 🕘 🔵 🚺	tem Sets - 3:32	- Item Set Finder (Bo	orgelt) (ITEMSET)
File Hilite Navigation View			
Table "default"	- Rows: 4305	Spec - Columns: 4	Properties Flow Variables
rabie acram		Spec columns.	Troperties Troit fundoies
[] ItemSet	ItemSetSize	I ItemSetSupport	D 💌 RelativeltemSetSupport%
["UV signature"]	1	265	80.06
["immune"]	1	168	50.755
["BRAF_Hotspot_Mutants"]	1	150	45.317
["immune", "UV signature"]	2	142	42.9
["BRAF_Hotspot_Mutants","UV signature"]	2	136	41.088
["keratin"]	1	100	30.211
["RAS_Hotspot_Mutants"]	1	92	27.795
["O"]	1	91	27.492
["hyper-methylated"]	1	91	27.492
["RAS_Hotspot_Mutants", "UV signature"]	2	86	25.982
["CpG island-methylated"]	1	85	25.68
["hypo-methylated"]	1	84	25.378
["MIR.type.2"]	1	83	25.076
["MIR.type.1"]	1	82	24.773
["MIR.type.3"]	1	81	24.471
["2"]	1	77	23.263
["BRAF_Hotspot_Mutants","immune"]	2	77	23.263
["hypo-methylated","UV signature"]	2	73	22.054
["0","UV signature"]	2	73	22.054
["MIR.type.4"]	1	72	21.752
["hyper-methylated", "UV signature"]	2	72	21.752
["normal-like"]	1	71	21.45
["BRAF_Hotspot_Mutants","immune","UV signature"]	3	71	21.45
["MIR.type.1", "UV signature"]	2	70	21.148
["MIR.type.2", "UV signature"]	2	66	19.94
["keratin", "UV signature"]	2	66	19.94
["CpG island-methylated", "UV signature"]	2	65	19.637
["MIR.type.3", "UV signature"]	2	64	19.335
["2"."UV signature"]	2	63	19.033

Figure 22. Item sets obtained using for the "Item Set Finder" node and ordered by their relative percentage of support.



UNIVERSIDAD DE GRANADA

abiertaugi





2. Extract and analyze the association rules

Association rules can be extracted with the "Association Rule Learner" node. Figure 13 shows the different parameters and configuration options available in the "Association Rule Learner" node, including the minimum number of items, minimum support, and minimum confidence levels. In this example, we set the minimum number of items to 3, a support of 0.015, and a confidence level of 80%.

Item column: [] Transaccion_SplitResultList 📀
n set settings
Minimum set size: 3
Minimum support: 0.015 🗘
Absolute number Percentage

Figure 33. The "Association Rule Learner" node options and parameters.

Figure 14 shows a scatter plot indicating how the rules generated are distributed according to two metrics. In this example, the plot illustrates the percentage of confidence and support, although other metrics such as lift could also be used by selecting which ones to display in the upper right menu.









Figure 44. Scatter plot of the rules generated by applying the "Association Rule Learner" node, displayed as a function of the support and confidence metrics.



CES





.

abiertaugr

REFERENCES

- Bakos, G. (2013). KNIME essentials. Packt Publishing Ltd.
- Blokdyk, G. (2019). KNIME a Complete Guide 2019 Edition. Emereo Pty Limited, 2019.
- McCormick, K. (2019). Introduction to Machine Learning with KNIME. linkedin.com.
- Silipo, R. (2016). Introduction to Data Analytics with KNIME: A Data Science Approach to Analytics. O'Reilly.
- Silipo, R., & Mazanetz, M. P. (2012). The KNIME cookbook. KNIME Press, Zürich, Switzerland.
- Strickland, J. (2016). Data Analytics Using Open-Source Tools. Lulu. com.





