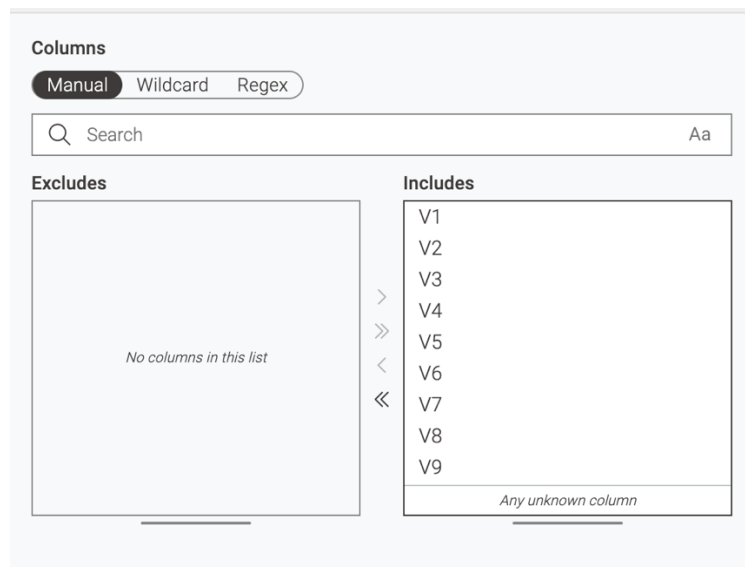
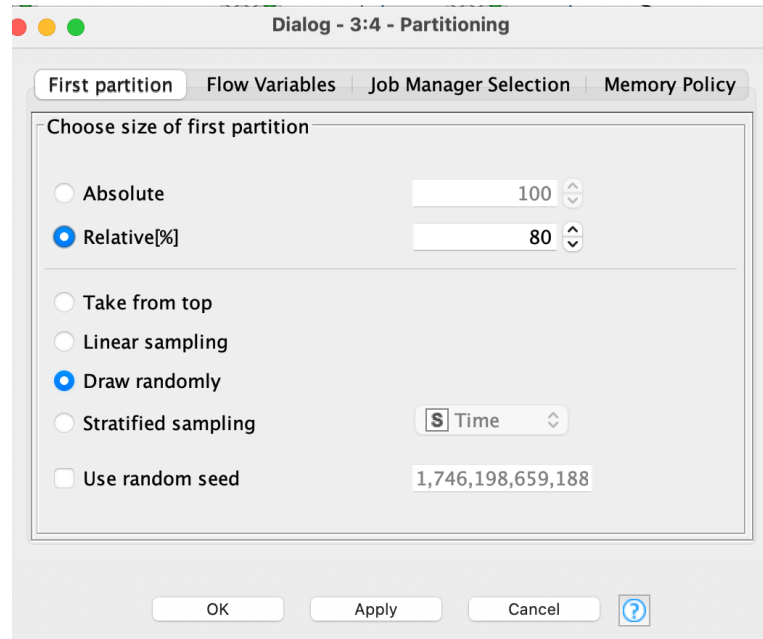


ابتدای کار از نود **number to string** استفاده میکنیم برای تبدیل داده های **numeric** به **string** به خاطر اینکه نود **decision tree learner** فقط بر روی داده هایی که تایپشون **string** هستند میتواند کار بکند



با توجه به تصویر تمام داده های **numeric** بودند که ما اذها را به **string** تبدیل کردیم



حالا با استفاده از نود **partition** داده هایمان را به **train** , **test** تقسیم میکنیم

در اینجا با انتخاب **80 % relative** ما 80 درصد از داده هایمان را به عنوان آموزش انتخاب کردیم

با انتخاب **draw randomly** هم گفتیم که بیاد و این 80 درصد از داده ها را رندوم انتخاب بکند

البته میتوانستیم **stratified sampling** را بر روی **class** و با یک **seed** مشخص هم انتخاب بکنیم که اگر این کار را میکردیم توزیع داده هایمان طبق **class** تقسیم میشد

General

Class column

Quality measure

Pruning method

☒ Reduced Error Pruning

Min number records per node

Number records to store for view

☒ Average split point

Number threads

☒ Skip nominal columns without domain information

Root split

☐ Force root split column

Root split column

Binary nominal splits

☐ Binary nominal splits

Max #nominal

☐ Filter invalid attribute values in child nodes

OK Apply Cancel ?

در نود decision tree learner هم تنظیمات را اینگونه انتخاب کردم

در قسمت class column ستونی که لیبل دار هست را انتخاب کردم و معیار تصمیم گیری هم ابتدا gini index انتخاب کردم

سپس این نود را به decision tree predictor وصل کردم که با توجه به الگوریتمی که بدست آوردیم لیبل کلاس را برای ما پیش بینی بکند

و بعد از آن برای اینکه بتوانیم accuracy مدل را بدست بیاورم به نود scorer وصل کردم .

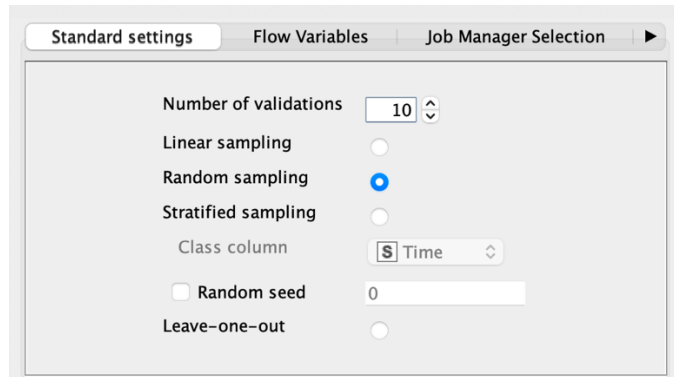
► 1: Confusion matrix ► 2: Accuracy statistics Flow Variables

Rows: 2 | Columns: 2

#	RowID	0 Number (integer)	1 Number (integer)
1	0	227451	0
2	1	394	0

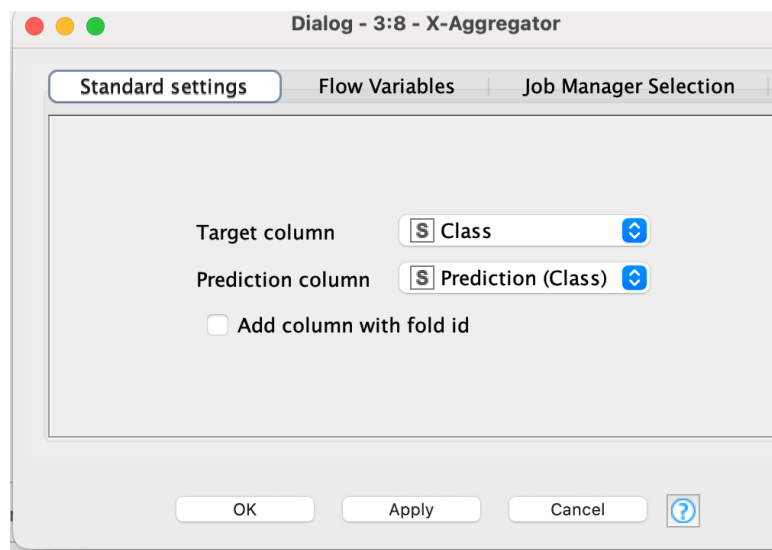
بخش دوم :

برای استفاده از cross validation هم دو نود x-partitioner , x-aggregator را اضافه کردم



در اینجا در بخش **number of validation** تعداد **validation** را انتخاب میکنیم که من به صورت دیفالت خوده ۱۰ را انتخاب کردم و گذاشتم داده هایم به صورت رندوم انتخاب بشوند.

بقیه ی مراحل هم مانند قبل هست نود **x-partitioner** را به **decision tree learner** میدهیم تا مدل ما را آموزش بده و بعد از اینکه آموزش داد نود **decision tree** را به **x-aggregator** وصل میکنیم .



درنود **x-aggregator** هم ستون هدف و ستون پیش بینی را انتخاب میکنیم .

بعد از این مورد برای ارزیابی این نود را به **scorer** وصل میکنیم

tows: 2 Columns: 2				
Table Statistics				
#	RowID	0 Number (integer)	1 Number (integer)	
1	0	284315	0	
2	1	492	0	

بخش سوم)

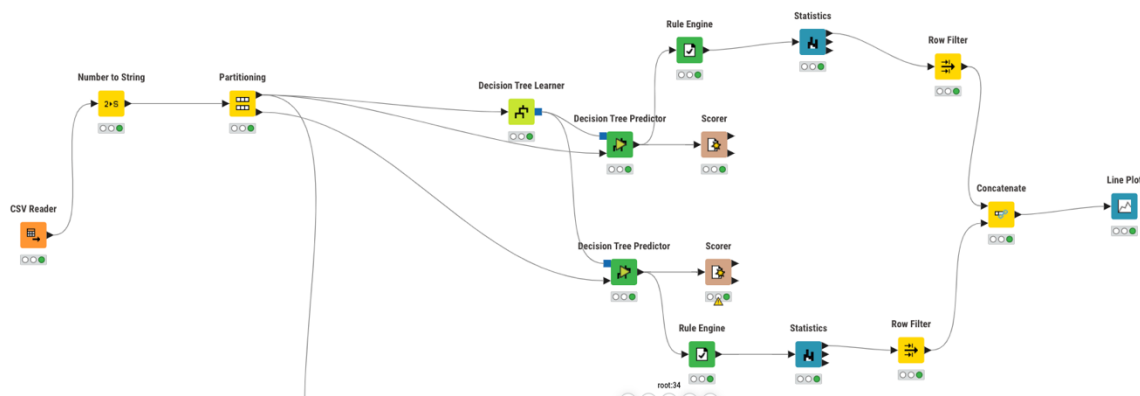
هرکدام از این معیار ها مزایا و معایب خاص خودش را دارد ولی اگر ما سرعت بیش تری در عملیات میخواهیم باید از معیار gini استفاده بکنیم چون سرعتش نسبت به entropy , gain ratio خیلی بیش تره و علتش هم بار عملیاتی کم ترش هست

ولی اگه دقت بیش تری میخواهیم معیار های entropy , gain ratio نسبت به gini خیلی بهتر عمل میکنند

یک مزیتی که gain ratio نسبت به gini , entropy دارد این است که دقت بهتر دارد و همچنین برای داده هایی که منحصر به فرد هستند مقاومت بالایی دارند در حقیقت Gain Ratio نرمالسازی شده است، بنابراین به ویژگی هایی که تعداد زیادی کلاس منحصر به فرد دارند، تنبیه می کند. در اینجا ما داده های منحصر به فرد نداشتیم که بخواهیم از gain ratio استفاده همچنین در این دیتاست دقت برای gini به اندازه ای خوب هست که به سراغ انترپی و gain ratio نرویم

سوال ۲ عملی)

برای این قسمت کل نود های ما به این صورت هست



بسیاری از نود ها و مراحل در قسمت قبلی توضیح داده شده

ولی برای نشانه دادن پچیدگی مراحل و overfitting از نود های rule engine , statistic , row filter , line plot استفاده کردم که در ادامه هرکدام از مراحل را توضیح میدم

در rule engine این دستورات وارد میکنیم برای اینکه مشخص بکنیم کدام رکورد ها غلط پیش بینی شده اند

برای داده های تست هم مانند train همین کار را میکنیم سپس با استفاده از line plot خطای داده های تست و ترین را باهم مقایسه میکنیم .

Rule Editor Flow Variables Job Manager Selection Memory Policy

Column List

- S ChestPainType
- S RestingBP
- S Cholesterol
- S FastingBS
- S RestingECG
- S MaxHR
- S ExerciseAngina
- S Oldpeak
- S ST_Slope
- S HeartDisease
- S Age
- D P (HeartDisease=0)
- D P (HeartDisease=1)
- S Prediction (HeartDisease)

Flow Variable List

- knime.workspace

Category

All

Function

- ? < ?
- ? <= ?
- ? = ?
- ? > ?
- ? >= ?
- ? AND ?
- ? IN ?
- ? LIKE ?
- ? MATCHES ?
- ? OR ?
- ? XOR ?
- FALSE

Description

Expression

```

1 // enter ordered set of rules, e.g.:
2 // $double column name$ > 5.0 => "large"
3 // $string column name$ LIKE "*blue*" => "small and blue"
4 // TRUE => "default outcome"
5 $HeartDisease$ = $Prediction (HeartDisease)$ => 0
6 TRUE => 1

```

Rows: 2 | Columns: 16

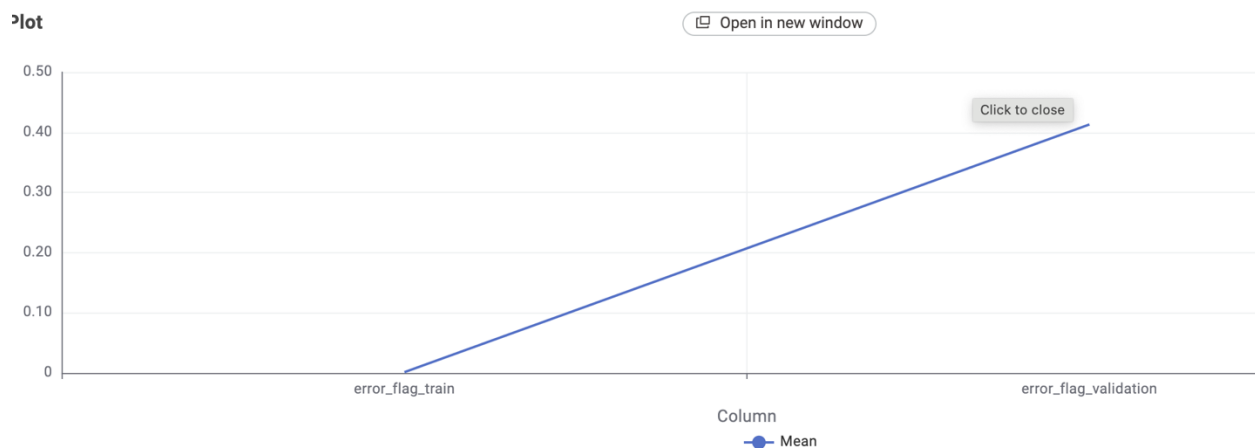
#	RowID	Column String	Min Number (dou...)	Max Number (dou...)	Mean Number (dou...)	Std. devia... Number (dou...)	Variance Number (dou...)	Skewness Number (dou...)	Kurtosis Number (dou...)	Overall su... Number (dou...)	No. missi... Number (inte...)	No. NaNs Number (inte...)	No. +ms Number (inte...)	No. -ms Number (inte...)	Median Number (dou...)	Row count Number (inte...)	Histogr SVG imag
1	Row0	error_flag_train	0	1	0.191	0.393	0.155	1.578	0.49	140	0	0	0	0	0	734	
2	Row1	error_flag_valid	0	1	0.168	0.375	0.141	1.786	1.203	31	0	0	0	0	0	184	

برای نشان دادن و محاسبه ی train , test مانند قبل دو نود decision tree predictor داریم که با استفاده از الگوریتم بدست آمده از decision tree learner و داده های تقسیم بندی شده از partitioner برای ما پیش بینی میکند

برای اینکه داده های ما overfit بشوند یکبار min number record per node را برابر ۱ قرار میدهیم خروجی مانند شکل پایین میشود همانطور که مشاهده میکنید generalization error ما نسبت به train error خیلی بیش تر است که به ما overfit شدن داده هایمان را نشان میدهد

Rows: 3 | Columns: 17

#	RowID	Column String	Min Number (dou...)	Max Number (dou...)	Mean Number (dou...)	Std. devia... Number (dou...)	Variance Number (dou...)	Skewness Number (dou...)	Kurtosis Number (dou...)	Overall su... Number (dou...)	No. missi... Number (inte...)	No. NaNs Number (inte...)	No. +ms Number (inte...)	No. -ms Number (inte...)	Median Number (dou...)	Row count Number (inte...)	Histogr SVG imag
1	P (Hea	P (HeartDisease	0	1	0.447	0.496	0.246	0.211	-1.958	328	0	0	0	0	0	734	
2	P (Hea	P (HeartDisease	0	1	0.553	0.496	0.246	-0.211	-1.958	406	0	0	0	0	0	734	
3	error_f	error_flag_train	0	1	0.001	0.037	0.001	27.092	734	1	0	0	0	0	0	734	



داده های train قبل از pruning با ۱ = min number recored

Rows: 2 | Columns: 2

#	RowID	0	1
1	0	328	0
2	1	1	405

Rows: 3 | Columns: 11

#	RowID	TruePositives	FalsePositives	TrueNegatives	FalseNegatives	Recall	Precision	Sensitivity	Specificity	F-measure	Accuracy	Cohen's kappa
1	0	328	1	405	0	1	0.997	1	0.998	0.998	0.998	0.997
2	1	405	0	328	1	0.998	1	0.998	1	0.999	0.999	0.997
3	Overall	0.998	0.998	0.998	0.998	0.998	0.998	0.998	0.998	0.999	0.999	0.997

داده های تست قبل از pruning با ۱ = min number record

Rows: 2 | Columns: 2

#	RowID	0	1
1	0	46	14
2	1	9	62

Rows: 3 | Columns: 11

#	RowID	TruePositives	FalsePositives	TrueNegatives	FalseNegatives	Recall	Precision	Sensitivity	Specificity	F-measure	Accuracy	Cohen's kappa
1	0	46	9	62	14	0.767	0.836	0.767	0.873	0.8	0.824	0.644
2	1	62	14	46	9	0.873	0.816	0.873	0.767	0.844	0.844	0.644
3	Overall	0.824	0.824	0.824	0.824	0.824	0.824	0.824	0.824	0.824	0.824	0.644

همانطور که از شکل و داده هایماز مشخص هست مدلا ما خیلی پیچیده شده است و برای داده های train ما دقت خیلی بالایی داریم ولی برای داده های تست دقت ما خیلی کم هست

بخش ۲)

حالا اگر در تنظیمات decision tree learner ما بیاییم و peruning method = mdl را انتخاب بکنیم خطای تست و ترین ما خیلی کمتر میشود

Dialog - 4:2 - Decision Tree Learner

Options PMMLSettings Flow Variables Job Manager Selection

General

Class column

Quality measure

Pruning method

☒ Reduced Error Pruning

Min number records per node

Number records to store for view

☒ Average split point

Number threads

☒ Skip nominal columns without domain information

Root split

☐ Force root split column

Root split column

Binary nominal splits

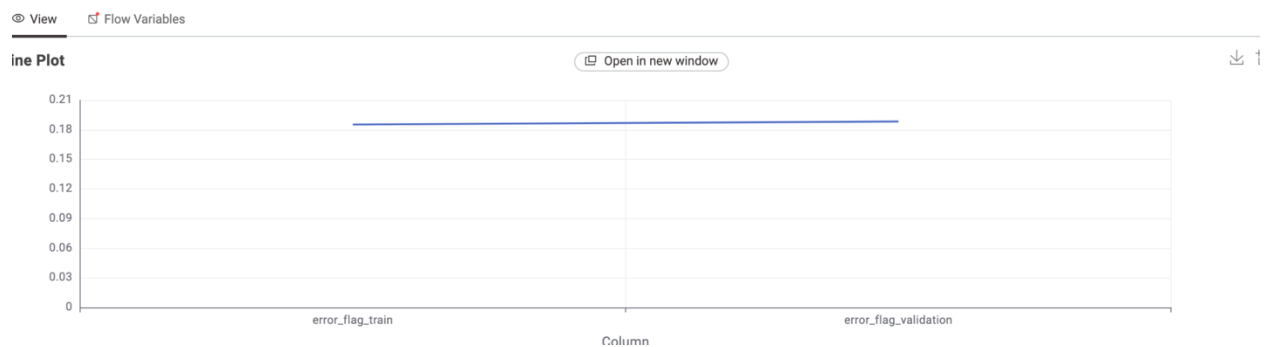
☐ Binary nominal splits

Max #nominal

☐ Filter invalid attribute values in child nodes

OK Apply Cancel ?

اگر تنظیمات ب این گونه تغییر بکند نتایج ما به صورت زیر است :



خطای داده های train بعد از pruning

Rows: 2 | Columns: 2

Table | Statistics

#	RowID	0 Number (integer)	1 Number (integer)
1	0	251	77
2	1	63	343

Rows: 3 | Columns: 11

Table | Statistics

#	RowID	TruePositives Number (integer)	FalsePositives Number (integer)	TrueNegatives Number (integer)	FalseNegatives Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's kappa Number (double)
1	0	251	63	343	77	0.765	0.799	0.765	0.845	0.782	0.77	0.77
2	1	343	77	251	63	0.845	0.817	0.845	0.765	0.831	0.831	0.831
3	Overall	0	0	0	0	0	0	0	Missing Value	0	0.809	0.613

داده های تست بعد از pruning:

Rows: 2 | Columns: 2

Table | Statistics

#	RowID	0 Number (integer)	1 Number (integer)
1	0	66	16
2	1	15	87

15

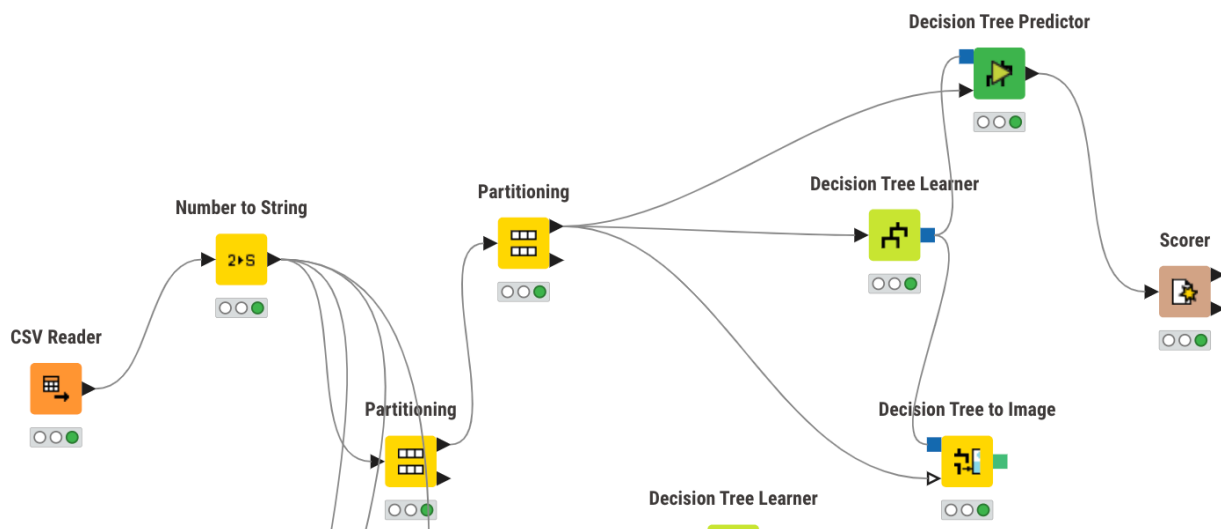
Rows: 3 | Columns: 11

Table | Statistics

#	RowID	TruePositives Number (integer)	FalsePositives Number (integer)	TrueNegatives Number (integer)	FalseNegatives Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's kappa Number (double)
1	0	66	15	87	16	0.805	0.815	0.805	0.853	0.81	0.81	0.81
2	1	87	16	66	15	0.853	0.845	0.853	0.805	0.849	0.849	0.849
3	Overall	0	0	0	0	0	0	0	0	0	0.832	0.659

همانطور که از تصاویر بالا مشخص هست pruning باعث میشود مدل ما overfit نشود و تعادل بین داده های train , test باشد

بخش ۳)



برای این بخش ما از دو partitioning استفاده کردیم در ابتدا ۸۰ درصد داده ها را برای train و ۲۰ درصد داده ها را برای test اختصاص می‌دهیم و سپس در partitioning دوم یک تقسیم بندی دوباره روی داده های train انجام می‌دهیم و از کم به زیاد درصد داده ها را زیاد می‌کنیم و بررسی می‌کنیم که خطای generalization ما چگونه است

وقتی از ۸۰ درصد داده های آموزش استفاده می‌کنیم خروجی نود scorer ب این صورت میشود :

#	RowID	TruePositives Number (integer)	FalsePositives Number (integer)	TrueNegatives Number (integer)	FalseNegatives Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's kappa Number (double)
1	0	222	21	304	40	0.847	0.914	0.847	0.935	0.879	0.909	0.788
2	1	304	40	222	21	0.935	0.884	0.935	0.847	0.909	0.909	0.788
3	Overall	0.909	0.909	0.909	0.909	0.909	0.909	0.909	0.909	0.909	0.896	0.788

وقتی از ۷۰ درصد از داده های آموزش استفاده می‌کنیم :

#	RowID	TruePositives Number (integer)	FalsePositives Number (integer)	TrueNegatives Number (integer)	FalseNegatives Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's kappa Number (double)
1	0	192	22	262	37	0.838	0.897	0.838	0.923	0.867	0.885	0.766
2	1	262	37	192	22	0.923	0.876	0.923	0.838	0.899	0.885	0.766
3	Overall	0.885	0.885	0.885	0.885	0.885	0.885	0.885	0.885	0.885	0.885	0.766

وقتی از ۵۰ درصد از داده های آموزش استفاده می‌کنیم :

برای این بخش هم مانند قبل ۲ قسمت `partitioner` داریم در `partitioner` اول ۸۰ درصد داده ها را برای آموزش داریم و در `partitioner` دوم میابیم از این ۸۰ درصد داده ای که جدا کردیم ۳۰ درصدش را برای `validation` و ۷۰ درصد را برای `train` انتخاب میکنیم و در ادامه سعی میکنیم با تغییر داده پارامترها بهترین مدل را که بیشترین دقت را دارد انتخاب کنیم

مثلا اولین تنظیم را اینگونه فرض میکنیم :

Dialog - 4:42 - Decision Tree Learner

Options

PMMLSettings

Flow Variables

Job Manager Selection

General

Class column

\$

HeartDisease

Quality measure

Gini index

Pruning method

No pruning

☒ Reduced Error Pruning

Min number records per node

15

Number records to store for view

10,000

☒ Average split point

Number threads

12

☒ Skip nominal columns without domain information

Root split

☐ Force root split column

Root split column

\$

ST_Slope

Binary nominal splits

☐ Binary nominal splits

Max #nominal

10

☐ Filter invalid attribute values in child nodes

OK

Apply

Cancel

?

[illegible]

مقدار دقت داده ی validation هم اینگونه است :

[illegible]

حال دوباره تنظیمات را تغییر میدهیم :

Dialog - 4:42 - Decision Tree Learner

Options PMMLSettings Flow Variables Job Manager Selection

General

Class column

Quality measure

Pruning method

☒ Reduced Error Pruning

Min number records per node

Number records to store for view

☒ Average split point

Number threads

☒ Skip nominal columns without domain information

Root split

☐ Force root split column

Root split column

Binary nominal splits

☐ Binary nominal splits

Max #nominal

☐ Filter invalid attribute values in child nodes

OK Apply Cancel ?

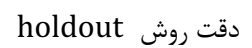
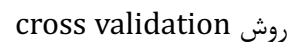
برای داده های تست مقدار دقت زیر را داریم :

[illegible]

برای داده های validation مقدار دقت زیر را داریم :

[illegible]

holdout روش



Rows: 3

Columns: 11

Table

Statistics

<input type="checkbox"/>	#	RowID	TruePositiv... Number (integer)	FalsePositi... Number (integer)	TrueNegati... Number (integer)	FalseNegat... Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's ka... Number (double)
<input type="checkbox"/>	1	0	68	19	83	14	0.829	0.782	0.829	0.814	0.805	?	?
<input type="checkbox"/>	2	1	83	14	68	19	0.814	0.856	0.814	0.829	0.834	?	?
<input type="checkbox"/>	3	Overall	?	?	?	?	?	?	?	?	?	0.821	0.639

دقت روش cross validation

Rows: 3

|

Columns: 11

Table

Statistics

<input type="checkbox"/>	#	RowID	TruePositiv... Number (integer)	FalsePositi... Number (integer)	TrueNegati... Number (integer)	FalseNegat... Number (integer)	Recall Number (double)	Precision Number (double)	Sensitivity Number (double)	Specificity Number (double)	F-measure Number (double)	Accuracy Number (double)	Cohen's ka... Number (double)
<input type="checkbox"/>	1	0	330	70	438	80	0.805	0.825	0.805	0.862	0.815	0.837	0.669
<input type="checkbox"/>	2	1	438	80	330	70	0.862	0.846	0.862	0.805	0.854	0.837	0.669
<input type="checkbox"/>	3	Overall	0.837	0.837	0.837	0.837	0.837	0.837	0.837	0.837	0.837	0.837	0.669

روش holdout اینگونه است که ما داده هایمان را به ۲ بخش تقسیم میکنیم مثلاً در اینجا من ۷۰ درصد داده ها را برای آموزش و ۳۰ درصد داده ها را برای تست قرار دادم حالا میام آموزش رو روی داده های آموزش انجام میدم و تست رو روی داده های تست انجام میدم

روش cross validation هم اینجوری هست که من داده هام رو به k بخش تقسیم میکنم مثلاً در اینجا به ۱۰ قسمت تقسیم کردم سپس مدلم به صورت چرخه ای بین بخش های مختلف آموزش داده می شود و سپس مدلم در بخش تست ارزیابی میشود

هرکدام از این روش ها مزایا و معایب خود را دارند، بنابراین انتخاب بهتر بستگی به شرایط پروژه و نیازهای خاص شما دارد.

اگر داده ها محدود یا کوچک باشند، Cross Validation معمولاً بهتر است، زیرا از تمام داده ها برای آموزش و ارزیابی استفاده می کند و دقت بیشتری در ارزیابی مدل ایجاد می کند.

اگر سرعت و سادگی مهم باشند، Holdout می تواند انتخاب بهتری باشد، زیرا سریع تر است و نیازی به اجرای چندین بار فرآیند آموزش ندارد.

به طور کلی، Cross Validation به دلیل دقت بالاتر و استفاده بهینه تر از داده ها ترجیح داده می شود، مگر در مواقعی که نیاز به سرعت و زمان کمتر دارید.