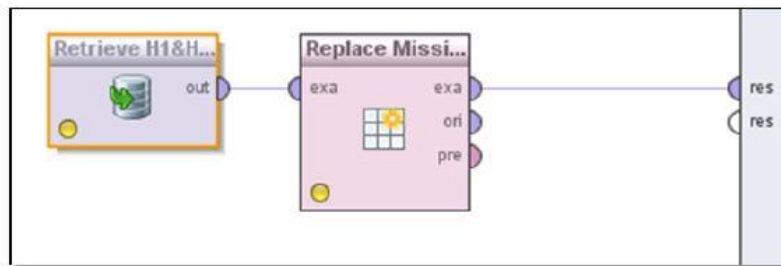


#### ۴-۳-۲- مقادیر مفقود شده:

تقریباً در تمام برنامه‌های کاربردی واقعی، می‌توان نمونه‌هایی را یافت که مقداری برای صفات خاصه آنها وجود ندارد. اجرای الگوریتم‌های داده‌کاوی بر روی داده‌های با مقادیر مفقوده زیاد، نتایج الگوریتم‌ها را می‌تواند نادرست کند. ساده‌ترین راه برای حل این مشکل صرف‌نظر از نمونه‌هایی است که برخی از صفات خاصه آنها دارای مقدار نیست. این در صورتی امکان‌پذیر است که تعداد رکوردهای داده‌ها بسیار زیاد باشد. در مواردی که تعداد ستون‌های یک متغیر کم باشند می‌توان به صورت دستی نمونه‌های ناقص را کامل کرد. که این راهکار برای پایگاه داده‌های حجیم امکان‌پذیر نیست.

در نرم افزار ریپدماینر ما می‌توانیم به صورت سیستمی و با چندین روش با داده‌های مفقوده رفتار کنیم. برای مثال می‌توان ابتدا از داده‌های ناقص صرف نظر نموده و نتیجه داده‌کاوی را مشاهده و پس از آن با جایگزین نمودن یک مقدار دوباره نتیجه حاصله را مشاهده و تحلیل کنیم. یکی از رایجترین عملگرهای ۱۱هایی که برای داده‌های مفقود شده می‌توان استفاده کرد Replace Missing Values می‌باشد. که ما می‌توانیم مقادیر مفقود شده را با مقادیر جدید به روش‌های مختلفی از جمله میانگین‌گیری، گذاشتن یک مقدار خاصی به جای مقادیر مفقود شده و ... جایگزین کنیم. ما در این تحقیق از این روش استفاده نموده ایم.



#### ۴-۳-۳- تغییر نوع داده:

شکل مناسب داده‌ها به عنوان ورودی الگوریتم‌های داده‌کاوی نقش به‌سزایی در این فرایند بازی می‌کنند و در مرحله آماده‌سازی داده‌ها این نقش پررنگ است. بسیاری از الگوریتم‌های داده‌کاوی برای نوع خاصی از صفات خاصه مناسب هستند و حتی برخی از آنها تنها بر روی نوع‌های مشخصی از صفات خاصه اجرا می‌شوند. برای مثال اکثر روش‌های رگرسیون بر روی داده‌های عددی اجرا می‌شوند. حال اگر در میان داده‌های ما داده‌ای از نوع اسمی باشد چگونه می‌توانیم از این روش‌ها برای تحلیل استفاده کنیم؟ عملگرهای متعددی در ریپدماینر وجود دارند که به ما کمک خواهد کرد تا این تبدیلات را انجام دهیم (شکل ۴-۲).

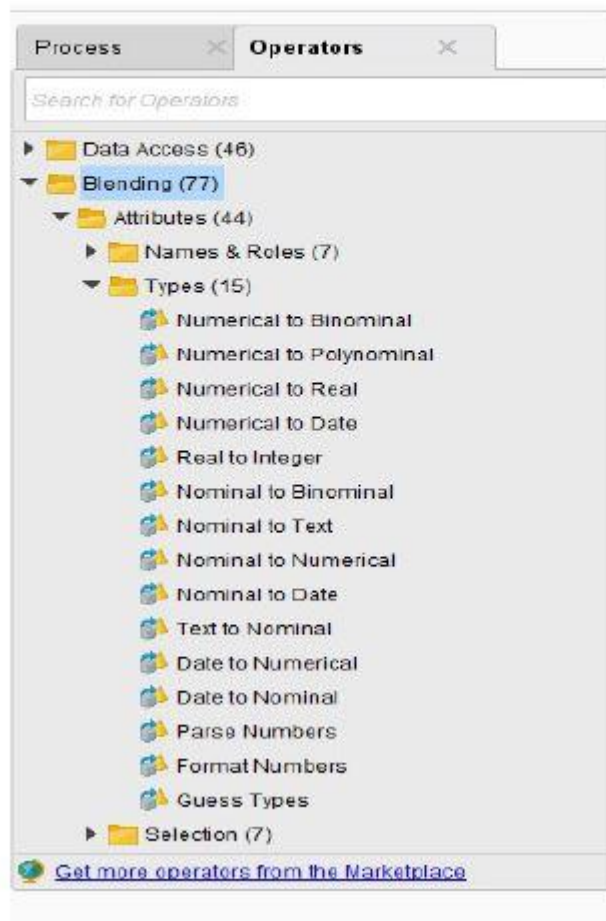
عملکرد بعضی از این عملگرها به شرح زیر می‌باشد.

صفات خاصه اسمی چند مقادیر را به نوع اسمی دو مقادیر ۱۲

صفات خاصه اسمی (چند مقداره) را به نوع عددی<sup>۱۳</sup>

صفات خاصه عددی را به اسمی (دومقداره)<sup>۱۴</sup>

صفات خاصه عددی را به اسمی (چند مقداره)<sup>۱۵</sup>



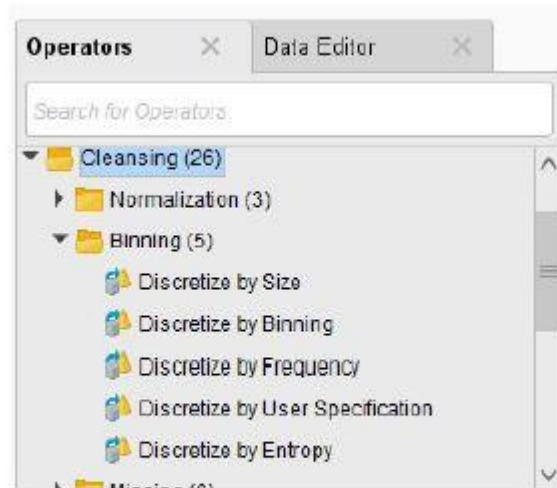
۴-۳-۴- نرمال سازی داده ها:

یکی دیگر از گونه‌های تغییر شکل داده‌ها نرمال‌سازی است. روش‌های نرمال‌سازی معمولاً برای مواردی که محاسبه فواصل بین نمونه‌ها در داده کاوی مطرح است مفید می‌باشد. در این روش مقادیر اندازه گیری شده به محدوده جدیدی نگاشت می‌شوند. این روش برای داده‌هایی که صفات خاصه آن واحد اندازه گیری متفاوتی دارند بسیار مناسب است.

#### ۴-۳-۵- گسسته سازی داده ها:

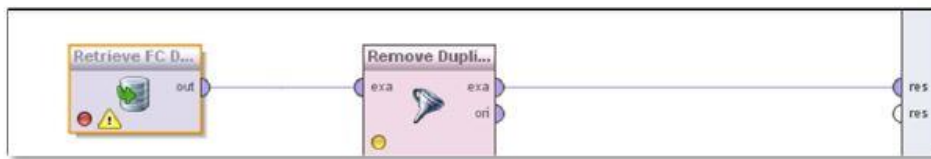
در واقع هدف گسسته‌سازی تبدیل داده‌های عددی به نوع اسمی است. با این عمل مقادیر یک صفت خاصه به چندین بازه تقسیم می‌شود. برای مثال فرض کنید ستونی از داده‌های شما نشان دهنده سن افراد است و نوع این صفت خاصه نیز عددی تعریف شده است. با عمل گسسته‌سازی شما می‌توانید این اعداد را به چندین بازه تقسیم کنید و برای هر بازه نامی را در نظر بگیرید برای مثال کلیه افراد زیر ۲۰ سال را در یک بازه و افراد بزرگتر از ۲۰ سال را در بازه دیگری قرار دهید بدین طریق شما محدوده‌ای از اعداد را در ۲ بازه قرار داده‌اید که می‌توانید با نام‌های جوان و میانسال آنرا جایگزین کنید. با این کار داده‌های عددی را به نوع اسمی تبدیل کرده‌اید. برخی از الگوریتم‌های داده کاوی با گسسته‌سازی کارا تر عمل می‌کنند. در ضمن این عمل برای نمایش قابل فهم داده‌ها نیز مفید است.

در نرم افزار ریپدماینر چندین روش برای انجام عمل گسسته‌سازی وجود دارد. این عملگرها در شکل زیر مشخص شده اند.



شکل (۴-۴) عملگر گسسته سازی داده ها در ریپدماینر

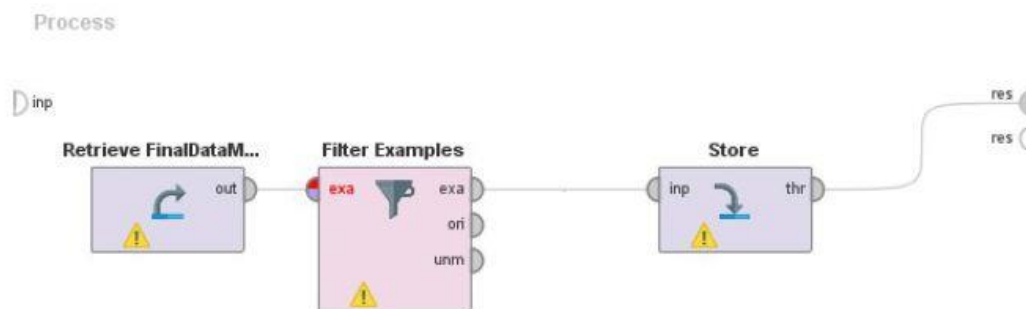
#### ۷-۳-۴- حذف نمونه های تکراری:



شکل (۸-۴) کاهش ابعاد داده ها (حذف سطرهای تکراری)

#### ۸-۳-۴- حذف برخی از نمونه ها<sup>۱۷</sup>:

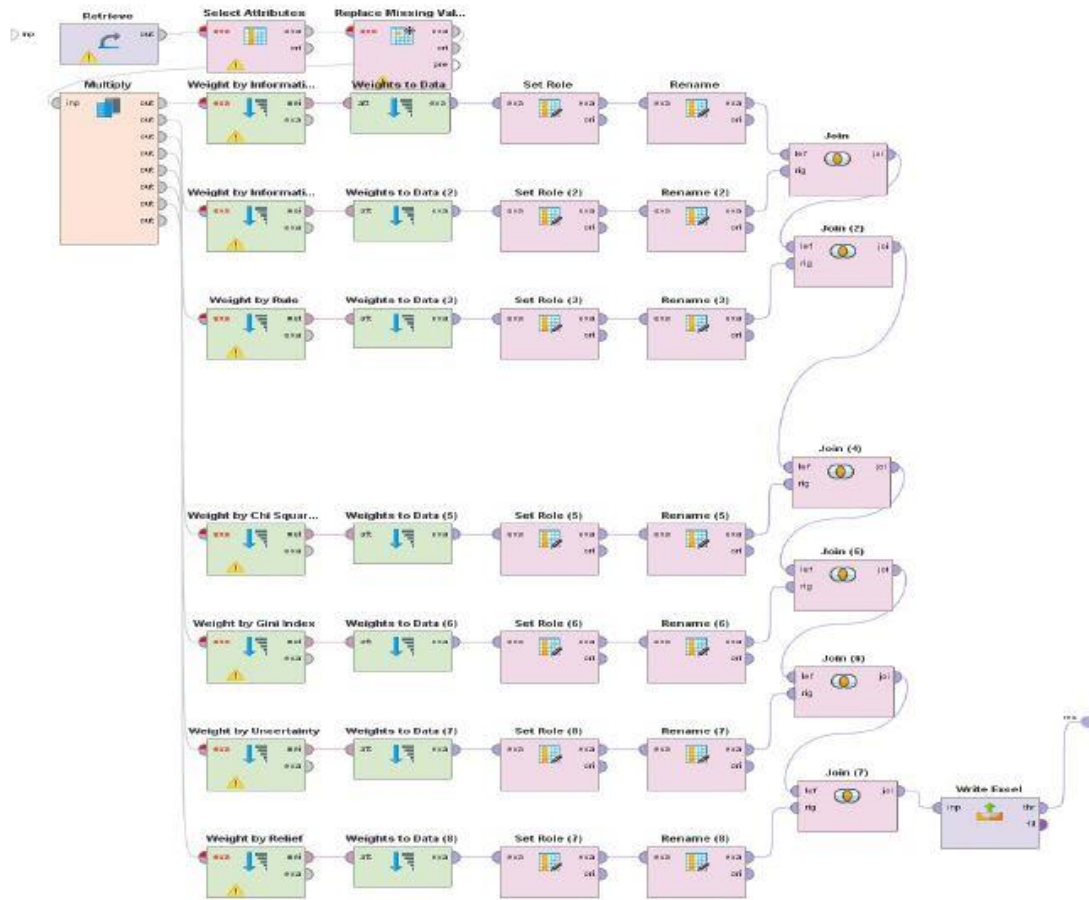
این عملگر در حقیقت وظیفه حذف برخی از نمونه ها را به عهده دارند.



شکل (۹-۴) حذف برخی از نمونه ها

#### ۴-۴- مدل وزن دهی<sup>۱۸</sup>:

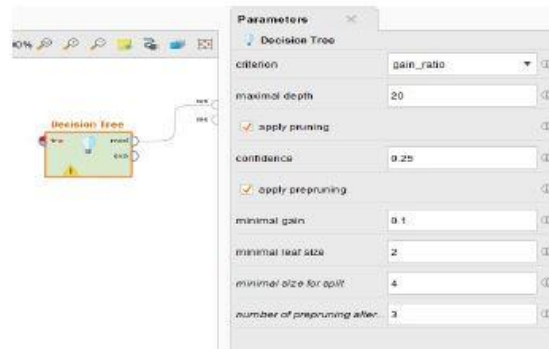
از طریق اعمال روش های وزن دهی جهت شناسایی مهمترین عوامل موثر بر پیشرفت تحصیلی استفاده شده است. در این تحقیق به دلیل ماهیت داده ها از ۷ عملگر وزن دهی در مدل وزن دهی استفاده شده است. این مدل در شکل زیر آمده است، همچنین نتایج اعمال روش های وزن دهی در بخش ۸-۴-۴ جدول ۲-۴ آمده است.

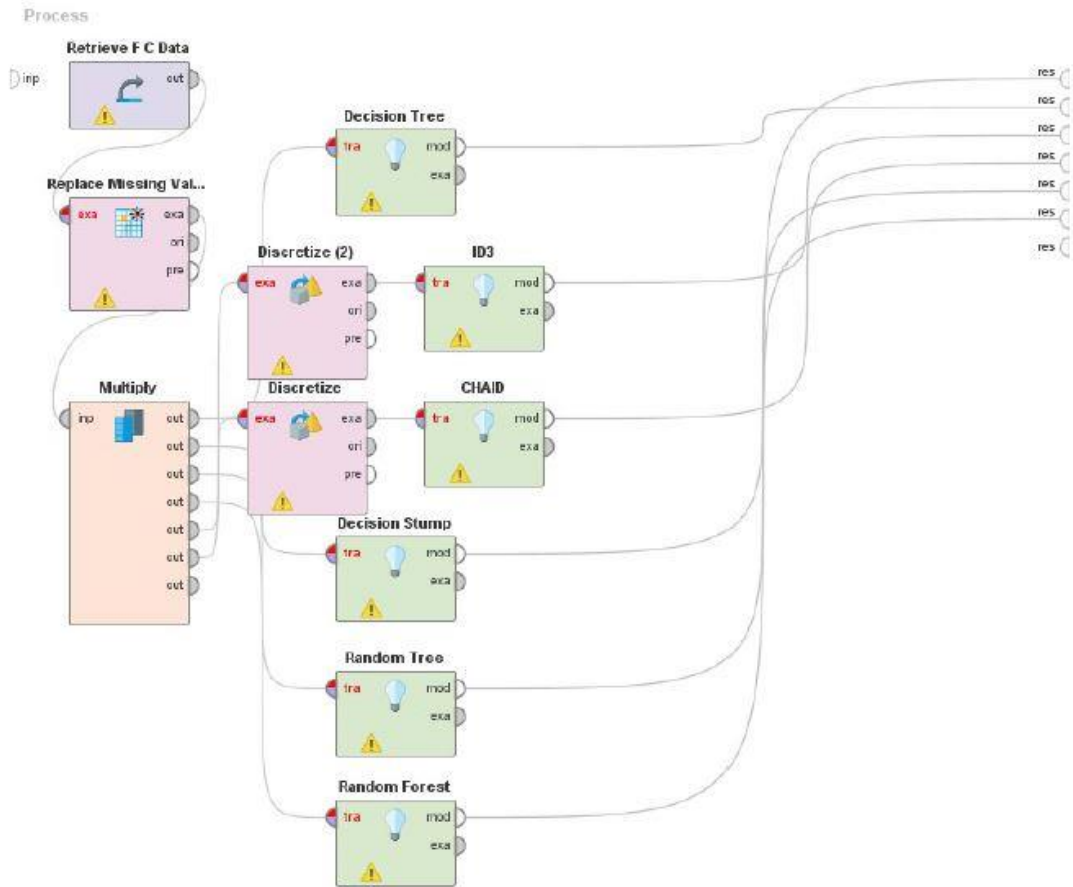


شکل (۴-۱۰) مدل وزن دهی

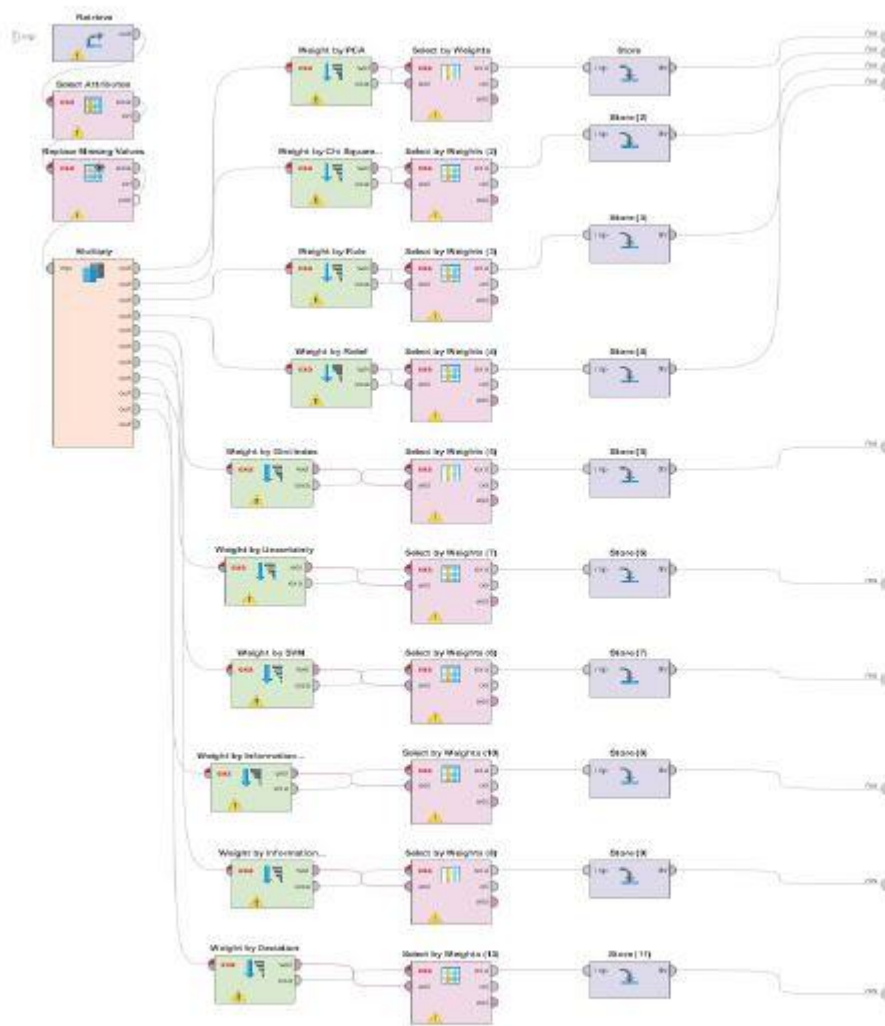
#### ۴-۵- استنتاج های درختی

روش یادگیری در این بخش روش Decision Tree است. و بقیه روش ها در واقع انشعابی از همین روش و یا انواع مقدماتی تر همین روش با قابلیت دست کاری کمتر هستند. برخی از درخت ها در این روش با استفاده از یک مقدار آستانه اقدام به هرس می کنند تا درخت از میزان خاصی که مورد نظر ما است وسیع تر نشود. این امر باعث می شود تا قدرت تفسیر نیز بالاتر برود.





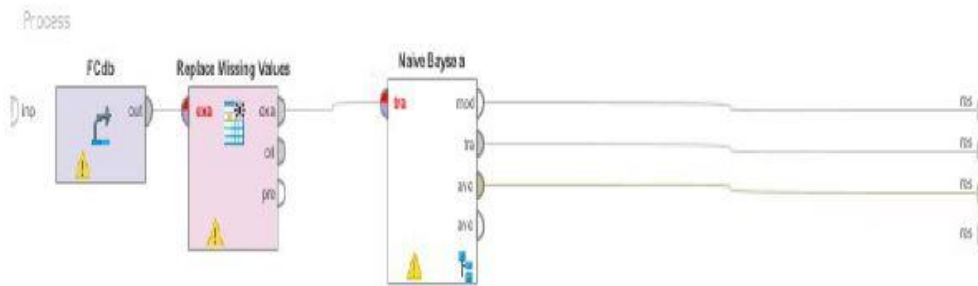
شکل (۴-۱۲) مدل رسم درخت تصمیم برای دسته بندی داده ها



شکل (۴-۱۶) مدل Attribute Selection

#### ۴-۶- مدل پیش بینی وضعیت تحصیلی

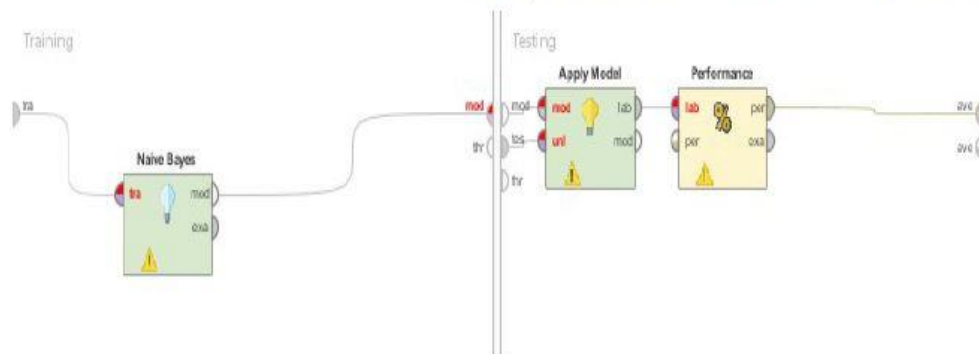
برای پیش بینی وضعیت های مختلف تحصیلی مدل زیر طراحی شد. این مدل که می توان انواع الگوریتم های مختلف برای پیش بینی وضعیت تحصیلی از آن استفاده کرد به راحتی وضعیت تحصیلی دانشجویان را پیش بینی می کند. در ادامه نتایج هر یک از الگوریتم های استفاده شده در این مدل عبارتند از الگوریتم های Decision Stump، Decision Tree، Naïve Bayes (Kernel) و یادگیری عمیق آمده است.



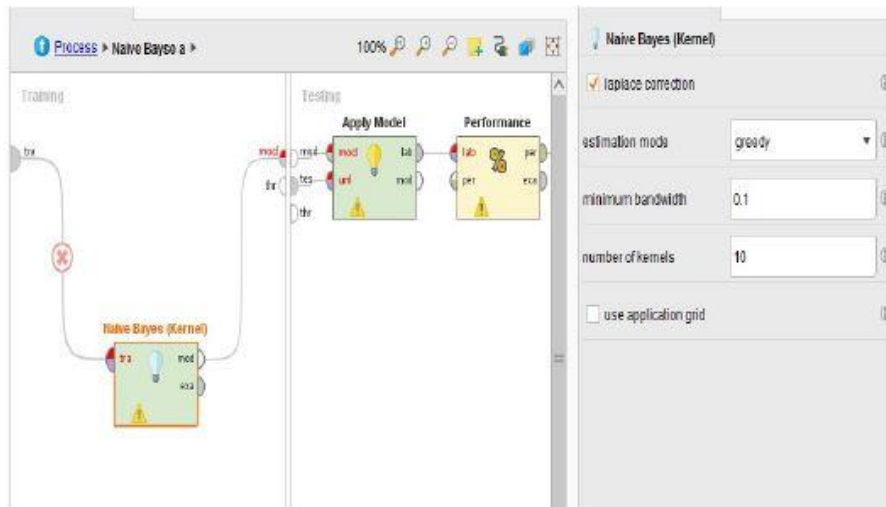
شکل (۴-۱۸) مدل پیش بینی نتایج

#### نتایج حاصل از اجرای الگوریتم Naïve Bayes:

این عملگر یک مدل طبقه بندی بیز ساده ایجاد می کند. بیز ساده براساس قانون احتمال شرطی بیز، برای طبقه بندی داده ها به کار برده می شود. از مزایای بیز ساده اجرای راحت و نتایج خوب برای بسیاری از کاربردها می باشد. این عملگر با استفاده از روش توزیع نرمال تخمینی و بر اساس نظریه Bayes، مدلی را برای دسته بندی رکوردهای آتی، در اختیار ما قرار می دهد. پارامتر Laplace تصحیح های لاپلاس را به منظور جلوگیری از تاثیر زیاد احتمالات صفر انجام می دهد.



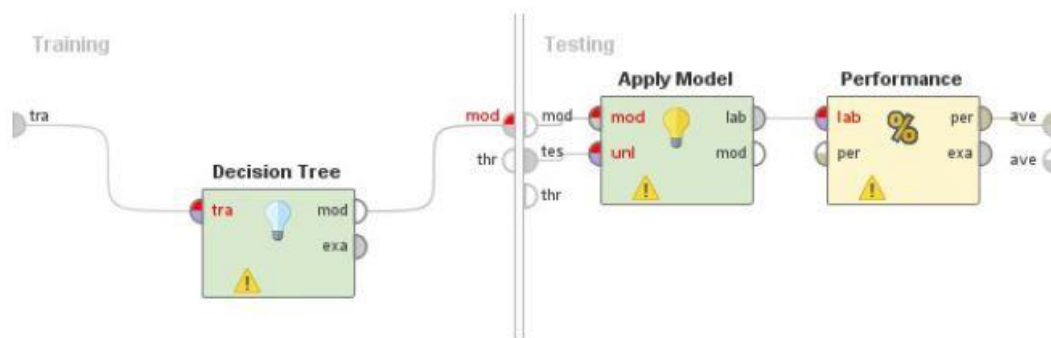
شکل (۴-۱۹) مدل پیش بینی نتایج توسط الگوریتم Naïve Bayes



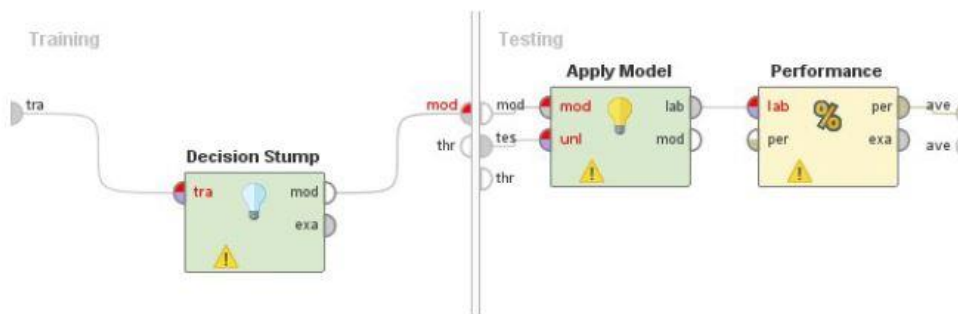
شکل (۴-۲۱) مدل پیش بینی نتایج توسط الگوریتم Naive Bayes (kernel)

#### ۴-۶-۳- نتایج حاصل از اجرای الگوریتم Decision Tree:

درخت تصمیم یکی از ابزارهای قوی و متداول برای دسته بندی و پیش بینی است. در ساختار درخت تصمیم، پیش بینی به دست آمده از درخت، در قالب یک سری قواعد توضیح داده می شود. در درخت تصمیم ضرورتی ندارد داده ها لزوماً از نوع عددی باشند داده ها می توانند متنی و اسمی هم باشند. در درخت تصمیم بالاترین گره در درخت گره ریشه است و گره های برگ، دسته ها یا توزیع دسته ها را نشان می دهند. ویژگی های مهم درخت تصمیم این است که استفاده از آن آسان است، درک مدل ایجاد شده توسط درخت تصمیم آسان است، در تقسیم بندی هیچ داده ای حذف نمی شود. توضیحات در مورد درخت تصمیم در قسمت استنتاج های درختی آمده است.



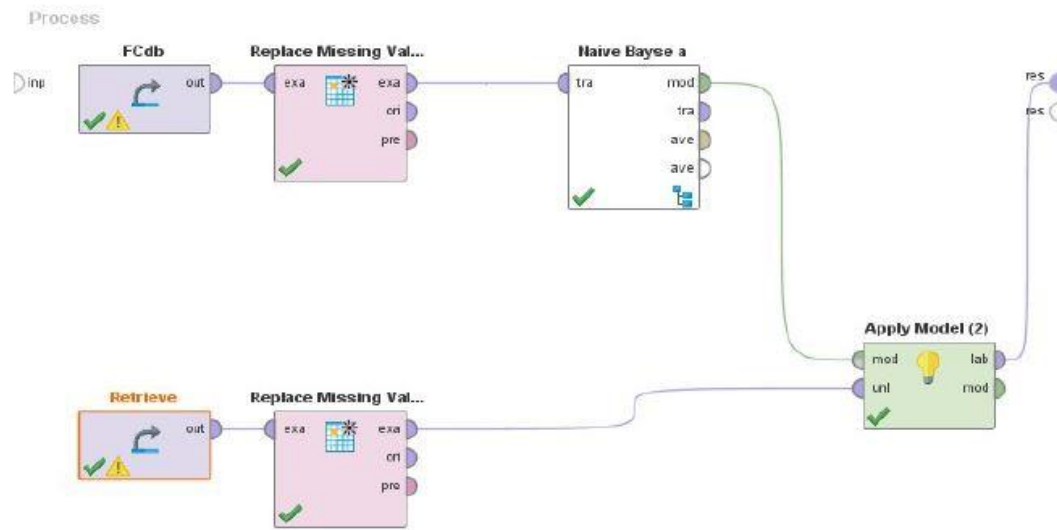
شکل (۴-۲۳) مدل پیش بینی نتایج توسط الگوریتم Decision Tree



شکل (۴-۲۵) مدل پیش بینی نتایج توسط الگوریتم Decision Stump

جدول (۴-۳) مقایسه نتایج مدل های پیش بینی وضعیت تحصیلی

دقت کلی	اخذ کاردانی		انصرافی		انتقالی		محروم از تحصیل		فارغ التحصیلی		الگوریتم
	صحت کلاس	فراخوانی کلاس	صحت کلاس	فراخوانی کلاس	صحت کلاس	فراخوانی کلاس	صحت کلاس	فراخوانی کلاس	صحت کلاس	فراخوانی کلاس	
84.57	۷.۷۳	20.93	۵۰.۵۰	40.48	۵۳.۵۷	67.98	۳۵.۴۶	51.79	۹۶.۸۹	91.45	Naive Bayes
79.54	۳.۹۱	53.94	۴۷.۸۳	40.81	۵۹.۵۲	69.01	۴۱.۲۸	44.39	۹۶.۶۶	85.55	Naive Bayes(Kemel)
91.18	۳۱.۸۲	8.14	۴۹.۷۰	54.03	۶۴.۳۴	68.86	۹۷.۱۶	45.92	۹۷.۱۵	99.16	Decision Tree
85.69		0		0	۰	0	۵۴.۰۰	57.55	۸۸.۰۰	99.6	Decision Stump
91.35	۱۰۰	2.23	۵۵.۶۱	40	۶۵.۵۱	71.64	۶۰.۳۰	55.61	۹۶.۹۰	99.34	Deep Learning



۴-۷- پیش بینی وضعیت تحصیلی دانشجویان جاری بر اساس مدل طراحی شده

Table view | Plot view

accuracy: 87.71% +/- 6.59% (mikro: 87.72%)

	true Fareghottahsil	true AkhzKardani	class precision
pred. Fareghottahsil	76	12	88.36%
pred. AkhzKardani	9	74	89.16%
class recall	89.41%	88.05%	

شکل (۴-۳۵) مدل پیش بینی پیشرفت وضعیت تحصیلی فارغ التحصیلی و اخذ کاردانی

accuracy: 97.02% +/- 1.61% (mikro: 97.02%)

	true Fareghottahsil	true Enseraf	class precision
pred. Fareghottahsil	608	25	96.05%
pred. Enseraf	12	595	98.02%
class recall	98.06%	95.07%	

شکل (۳۶-۴) مدل پیش بینی پیشرفت وضعیت تحصیلی فارغ التحصیلی و انصرافی

Table View  Plot View

accuracy: 96.34% +/- 1.47% (mikro: 96.34%)

	true Fareghottahsil	true Enteghali	class precision
pred. Fareghottahsil	657	24	96.48%
pred. Enteghali	26	860	98.21%
class recall	96.19%	96.49%	

شکل (۳۷-۴) مدل پیش بینی پیشرفت وضعیت تحصیلی فارغ التحصیلی و انتقالی