

1 تحلیل و درك مفاهيم داده

1.1 انتزاع و تعميم داده ها

1

1. معدل كل (GPA): عددی پیوسته (Interval)
2. حضور در کلاس‌ها: عددی گسسته (Ratio)
3. میزان استفاده از منابع آموزشی مجازی و حضوری: (Ratio)
4. رشته تحصیلی: داده اسمی (Nominal)
5. سطح تحصیلات والدین: داده ترتیبی (Ordinal)
6. وضعیت اقتصادی خانواده: عددی گسسته (Interval/Ratio)
7. میزان تعامل با اساتید و دانشجویان: داده ترتیبی (Ordinal)
8. وجود مشکلات روانی یا استرس: (Nominal)

برای تحلیل احتمال ترک تحصیل دانشجویان نیازمند ویژگی‌های مختلفی هستیم که بتوان تخمین دقیق‌تری از ترک تحصیل بدست آوریم پس لازم است که عوامل تاثیرگذار در این ماجرا را بدست آوریم. به عنوان مثال معدل و میزان حضور شخص در کلاس می‌تواند دید خوبی از وضعیت تحصیلی و تمایل به ادامه تحصیل شخص بدهد هر چقدر یک دانشجو معدل پایین‌تر داشته باشد و تمایلی به شرکت در کلاس‌های درس نداشته باشد احتمال ترک تحصیل آن می‌تواند افزایش یابد.

از طرفی حمایت‌های تحصیلی از جمله دسترسی به منابع آموزشی و یا وضعیت اقتصادی خانواده در امر ادامه تحصیل بسیار اثرگذار می‌باشد. پس اگر دانشجویی شرایط مالی مطلوبی نداشته باشد و به منابع آموزشی دسترسی نداشته باشد احتمال ترک تحصیل کردن آن حتی در صورتی که معدل بالایی داشته باشد افزایش می‌یابد. همچنین رشته تحصیلی شخص که بیانگر بازار کار آینده دانشجو می‌باشد هم می‌تواند عامل اثرگذاری باشد.

که هر کدام از این داده‌ها بر ماهیتشان با نوع مختلفی از ویژگی‌ها تعریف می‌شوند، مثلاً معدل دانشجو یک داده پیوسته است و یا حضور در کلاس باید به صورت اعداد گسسته تعریف شود. و یا مثلاً برای تعامل با دانشجویان و اساتید باید از داده‌های ترتیبی استفاده کرد (مثلاً بد، خوب، خیلی خوب).

2

خبر همه ویژگی‌ها از اهمیت یکسانی برخوردار نیستند، به عنوان مثال اگر در میان ویژگی‌ها نام و کدملی دانشجویان قرار داشته باشد این موارد تاثیری در پیش‌بینی ما نخواهد داشت. و یا اینکه دانشگاه فرد نسبت به وضعیت مالی و نمرات دانشجو اهمیت کمتری داشته باشد.

برای بررسی اینکه کدام ویژگی اهمیت بیشتری دارد هم می‌تواند از دانش خودمان استفاده کنیم مثلاً مواردی که بالاتر گفته شده است و هم می‌توان از تکنیک‌هایی همون تحلیل همبستگی استفاده کرد در این روش میزان ارتباط بین ویژگی‌ها و تغییر هدف (مثل ترک تحصیل) با استفاده از یک ضریب همبستگی یا heatmap بررسی می‌کنیم. و یا می‌توان از روش‌های انتخاب ویژگی مانند Mutual استفاده کرد.

2.1 اثر نوع ویژگی بر تحلیل داده ها

1

خیر این دو ویژگی را نمی‌توان به یک صورت در تحلیل های آماری در نظر گرفت زیرا نوع داده ها متفاوت است، مثلاً درآمد سالانه یک ویژگی عددی پیوسته است که می‌توان بر روی آن عملیات ریاضی همچون میانگین و واریانس و یا روش های تحلیل همبستگی اعمال کرد و حتی می‌توان داده های مختلف در درآمد سالانه را باهمدیگر مقایسه کرد.

اما از جهتی کدپستی یک ویژگی گسسته میباشد که صرفاً نشان دهنده یک خانه بوده و مقایسه آن با دیگر داده‌های کد پستی اطلاعاتی در اختیار ما قرار نمی‌دهد حتی می‌تواند ما را به گمراهی بکشاند از طرفی بر روی این نمونه داده های نمی‌توان عملیات ریاضی همچون واریانس، میانگین و مد اعمال کرد چرا که نتیجه یک عدد می‌شود که هیچ دید آماری به ما نمی‌دهد. صرفاً یک عدد بی معنی تولید خواهد شد.

2

برای نرمال سازی داده هایی مثل درآمد سالانه می‌توان از روش هایی همچون Min-Max Scaling و Z-score Normalization استفاده کرد. منتها چون در این روش های محاسبات ریاضی همچون میانگین صورت می‌گیرد در سوال قبل گفتیم که این نوع روش برای داده هایی مثل کدپستی جوابگو نیست چرا که میانگین گرفتن و یا تفریق و تقسیم چند داده مختلف از نوع کد پستی معنای خاصی به ما نمیدهد و حتی ممکن است در این بین اعداد پیوسته نامفهومی تولید شود.

اما برای اینکه داده ها در بازه بین صفر تا یک واقع شوند می‌توان به هر کدام از کد پستی ها یک عدد در بازه صفرو یک یا به صورت رندم و یا به صورت ترتیبی نسبت داد به عنوان مثال ، کد پستی اول مثلاً مقدار صفر را به آن نسبت دهیم کد پستی دوم مقدار 0.001 و کد پستی دوم مقدار 0.002 به همین شکل تا آخر، به عبارتی برای نرمال کردن اینگونه داده های می‌توان از روش های label گذاری استفاده کنیم. که همان طور که مشخص است از روش های محاسباتی و آماری برای نرمال کردن اینگونه داده ها استفاده نشده است.

3.1 نقد معیار توصیفی داده ها

1

وقتی میانگین و میانه برابر باشند اما واریانس بسیار بالا باشد، نشان دهنده یک توزیع متقارن با پراکندگی زیاد است. این به این معنی است که داده ها در اطراف مقدار میانگین توزیع شده اند اما فاصله زیادی از میانگین دارند. این یعنی ممکن است که داده ها از یک توزیع پهن مانند توزیع یکنواخت گسترده یا توزیع نرمال با انحراف معیار بالا باشند.

2

خیر میانگین در این حالت معیار مناسبی برای ما نیست. دلایلش هم این است که داده ها به دلیل واریانس زیاد از پراکندگی شدیدی برخوردار هستند . پس خود میانگین به تنهایی تصویری دقیق از تمرکز داده ها ارائه نمی‌دهد.

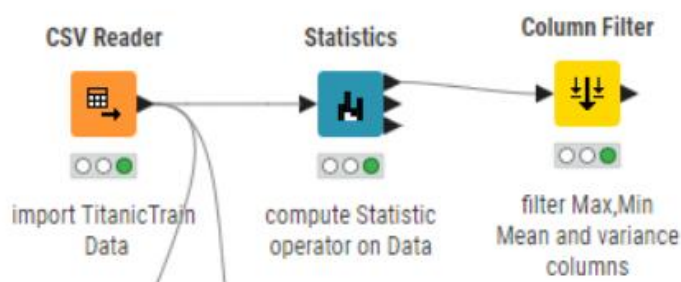
در این حالت می‌توان علاوه بر میانگین از انحراف معیار هم استفاده کرد و یا میتوان از چاک هم استفاده کرد چون داده‌ها را به باه‌هایی تقسیم می‌کند و می‌تواند دید آماری بهتری به ما بدهد. از طرفی برای نشان دادن توزیع داده‌ها میتوان از نمودارهایی همچون هیستوگرام نیز بهره برد.

گزارش بخش عملی

2 تحلیل مجموعه داده تایتانیک

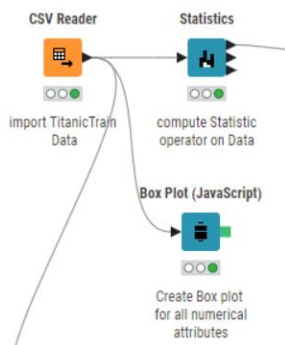
2.1 بصری سازی

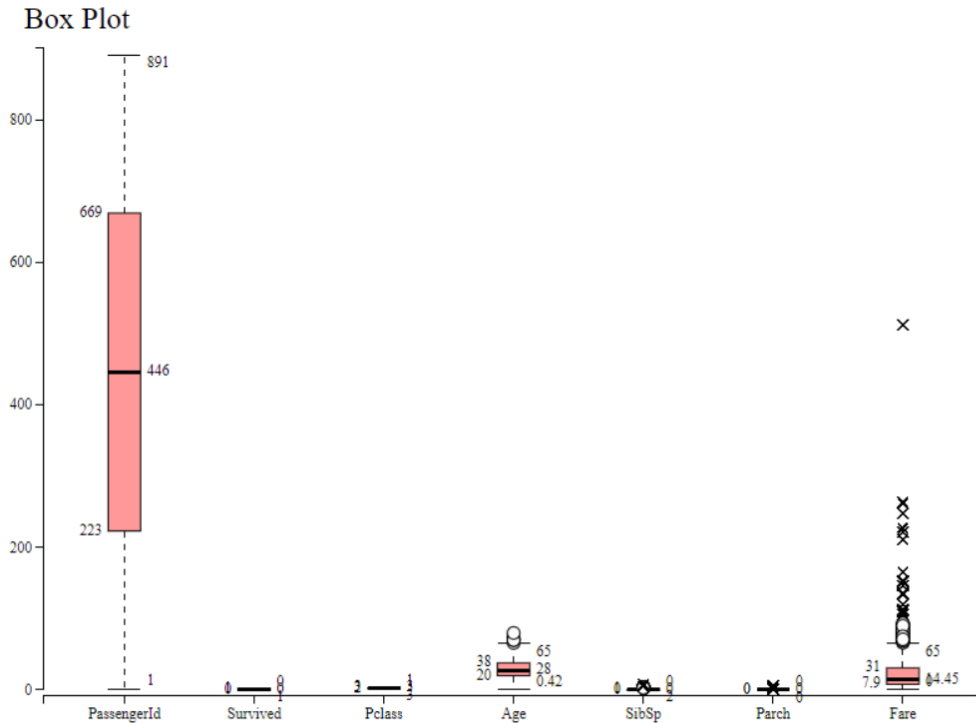
1



#	RowID	Min Number (double)	Max Number (double)	Mean Number (double)	Variance Number (double)
1	Passer	1	891	446	66,231
2	Survive	0	1	0.384	0.237
3	Pclass	1	3	2.309	0.699
4	Age	0.42	80	29.699	211.019
5	SibSp	0	8	0.523	1.216
6	Parch	0	6	0.382	0.65
7	Fare	0	512.329	32.204	2,469.437

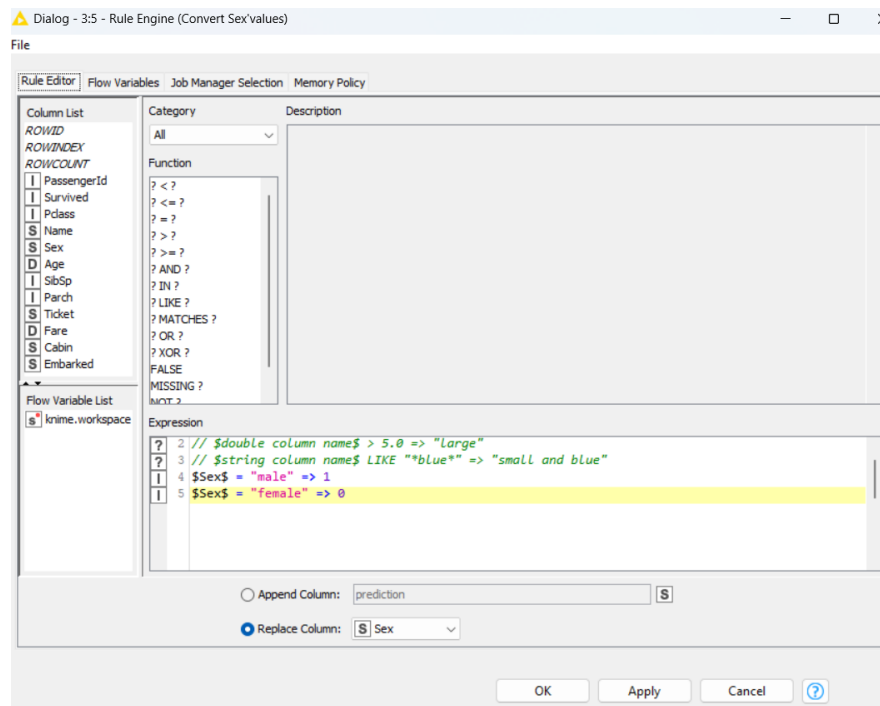
2

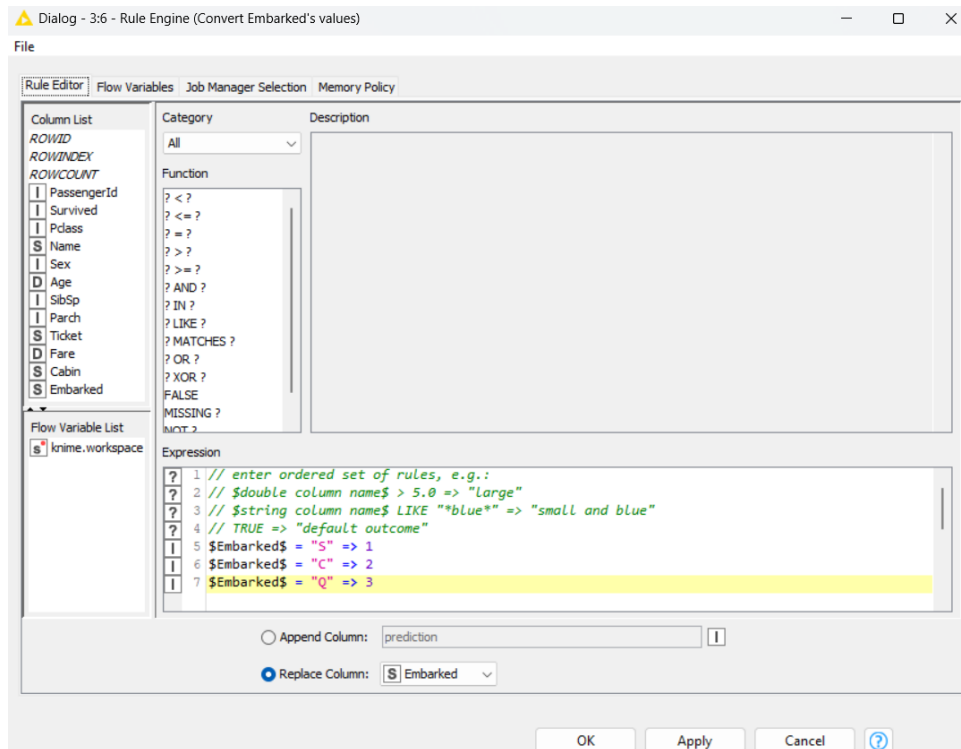




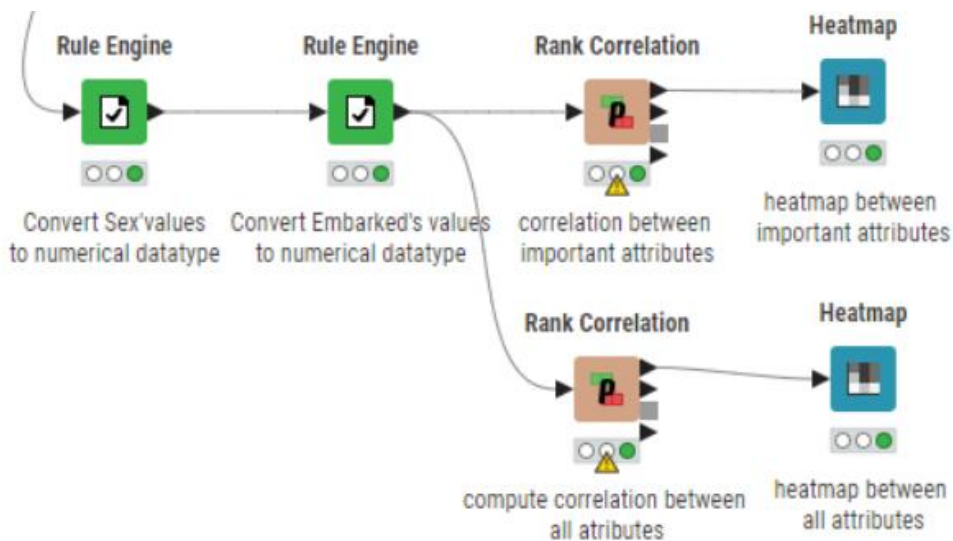
3

برای این سوال ابتدا ویژگی های Sex و Embarked را از حالت string با استفاده از Rule Engine به integer تبدیل میکنیم تا ارزیابی بهتری داشته باشیم .





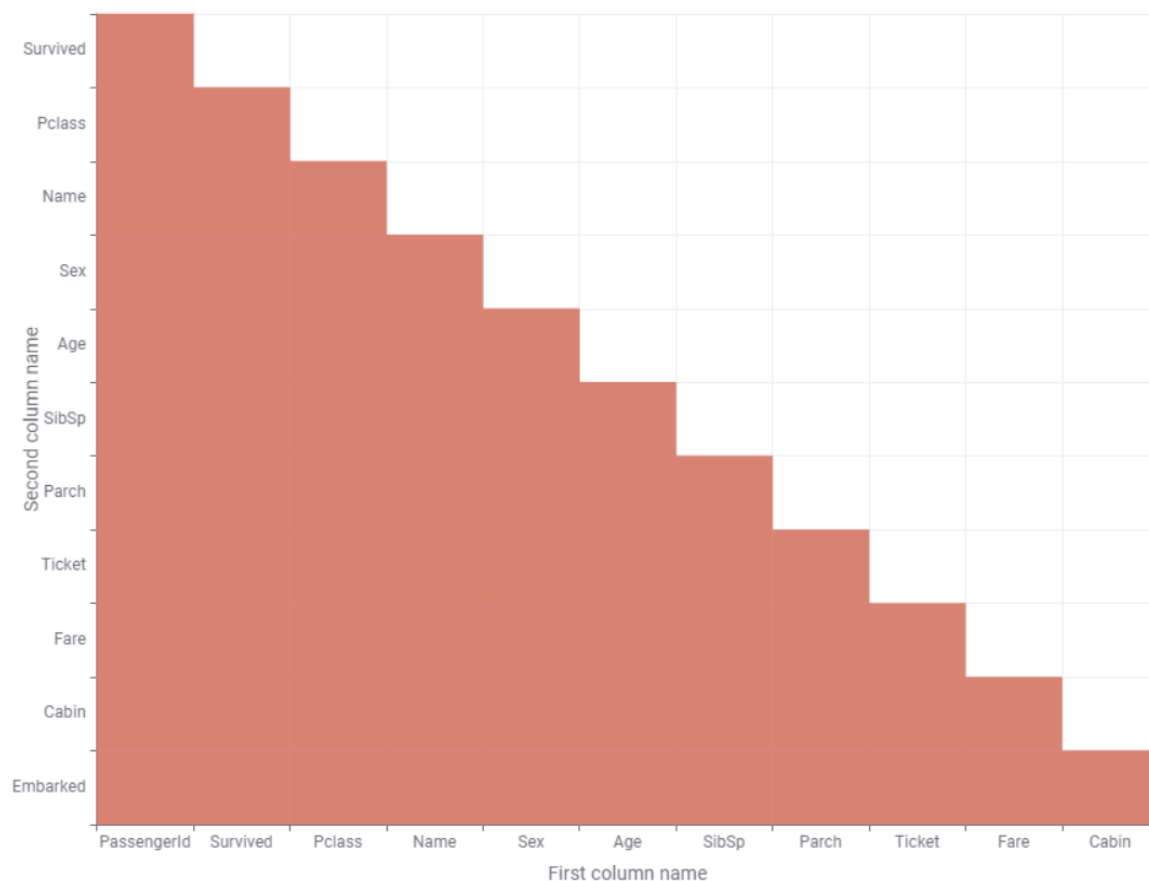
سپس این سوال را در دو حالت بررسی کرده ایم:



حالت اول اینکه ضریب همبستگی همه ویژگی ها حتی اسم و ایدی مسافران را هم در نظر گرفته ایم و پس از به دست آوردن ضریب همبستگی به نمودار heatmap داده ایم تا متوجه شویم بین کدام دو ویژگی بیشترین همبستگی را داریم.

با توجه به نمودار heatmap بدست آمده اگر تمامی attribute ها را لحاظ کنیم ، بیشترین همبستگی میان دو attribute، PassengerID و Survived می باشد.

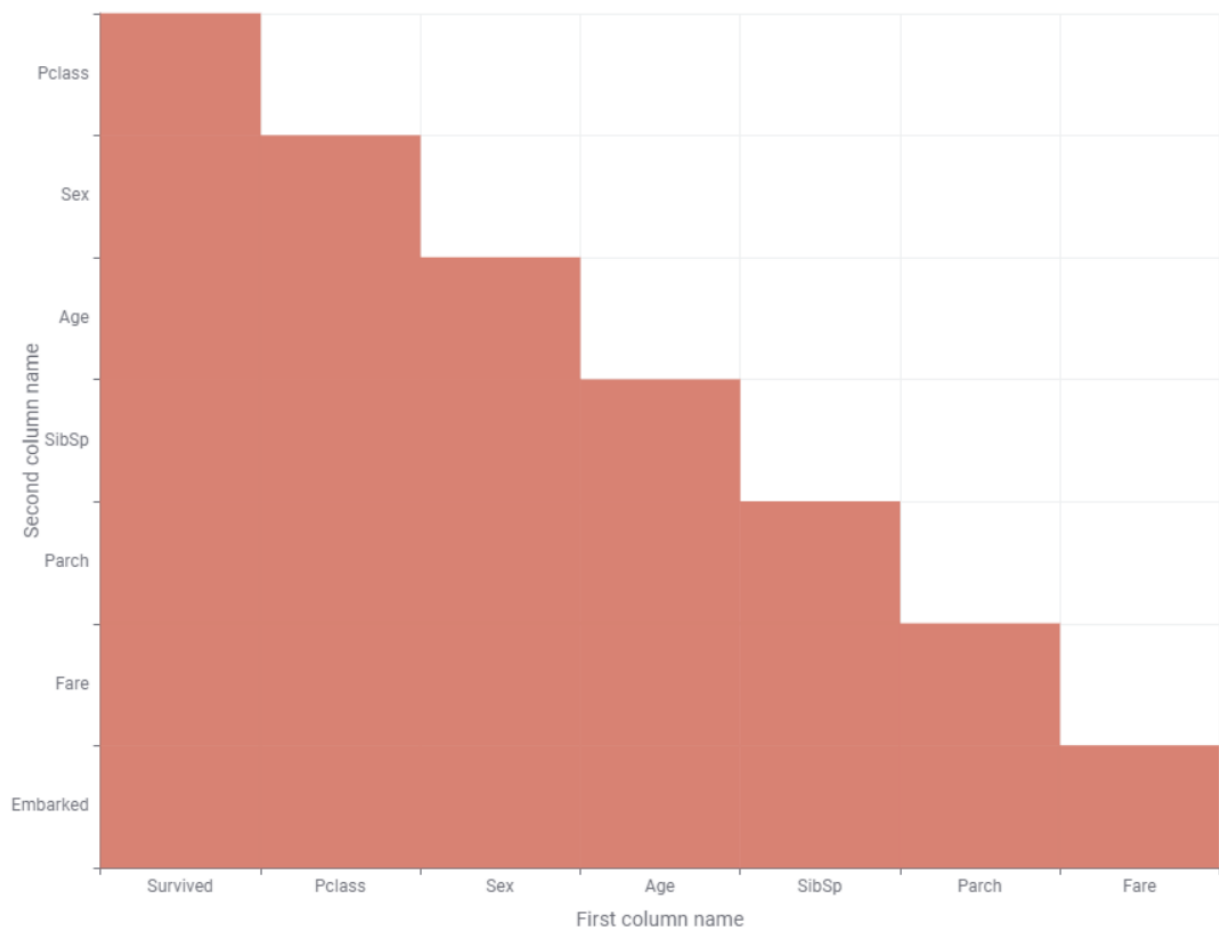
Heatmap



در حالت دوم صرفاً ویژگی‌های مهم را در نظر گرفتیم چرا که داده‌هایی مثل (Name, PassengerID, Ticket, Cabin) اطلاعات زیادی درباره مسافران در اختیار ما نمی‌گذارند و طبیعتاً برای هر مسافر می‌تواند به مقدار متناهی به فرد داشته باشد، داده‌کابین هم به این دلیل در نظر نمی‌گیریم چون مقدارهای تعریف نشده زیادی در فیلدهای این attribute هست.

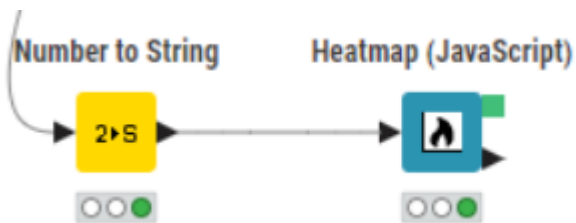
با توجه به نمودار بدست آمده برای حالت دوم بیشترین همبستگی میان دو ویژگی Survived و Pclass وجود دارد.

Heatmap



4

برای اینکه بتوانیم سطر ها را به عنوان PassengerID در نظر بگیریم باید ابتدا این ویژگی را به حالت string تبدیل کنیم. سپس به heatmap بدهیم.



برای نشان دادن مقادیر missing در تنظیمات خود heatmap مشخص شده که این مقادیر باید با رنگ سیاه نشان داده شوند.

Dialog - 3:12 - Heatmap (JavaScript)

File

Options View Configuration Interactivity Flow Variables Job Manager Selection Memory Policy

General

Width (px): 800

Height (px): 600

☒ Create image at output

☒ Show warnings in view

☒ Resize view to fill window

☒ Display fullscreen button

Display

Chart title:

Chart subtitle:

☐ Show tool tips

Gradient colors

☐ Use discrete gradient

Number of colors (odd): 3

Select gradient colors:

Change...

Change...

Change...

Missing value color

Select color for missing values: Change...

Out of range value colors

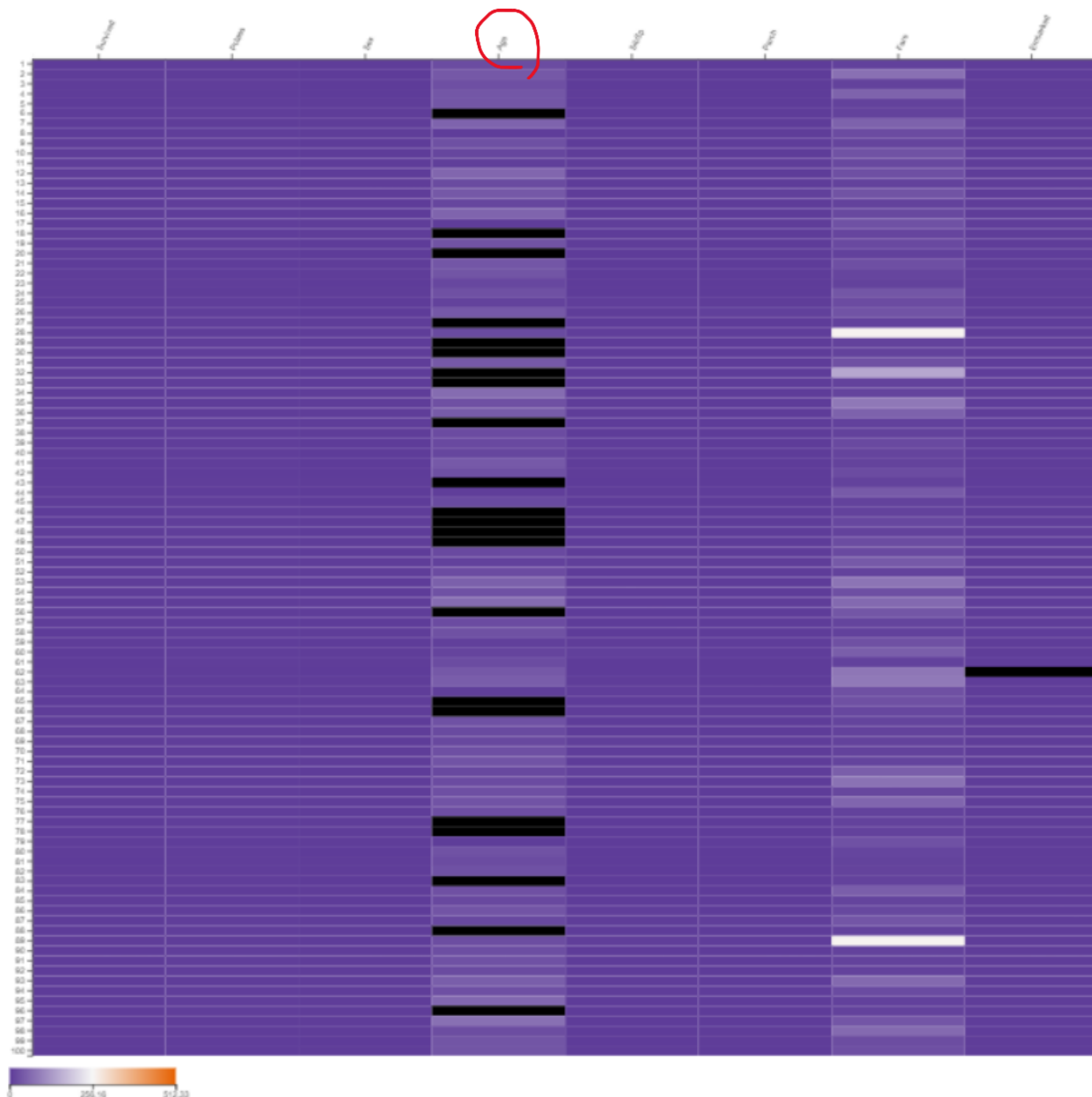
Select a color for values greater than the range maximum: Change...

Select a color for values less than the range minimum: Change...

OK Apply Cancel ?

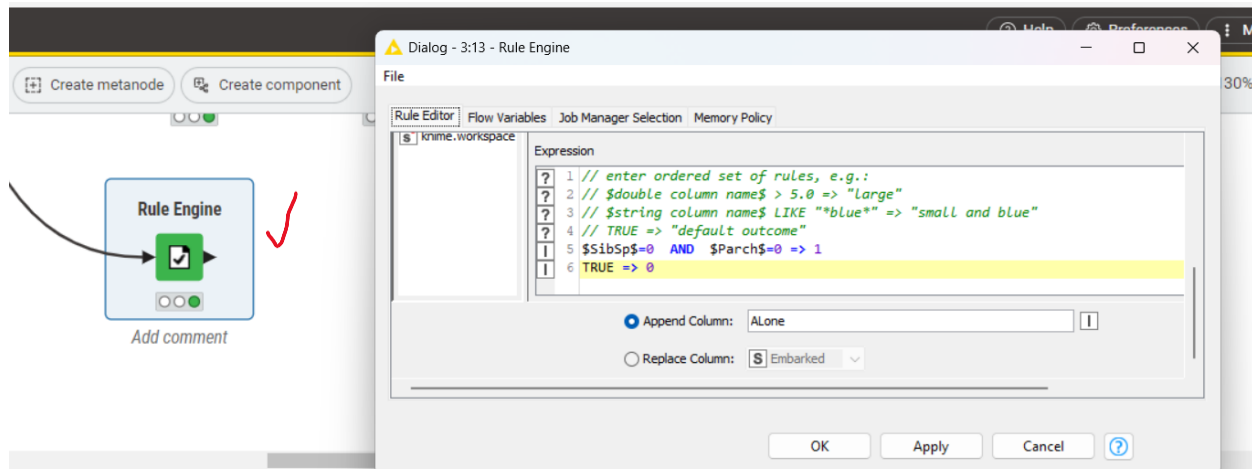
با توجه به خروجی بدست آمده فیلد هایی که مشکلی شده اند، missing value هستند و تعداد این مقادیر missing در ستون Age از باقی ستون ها بیشتر می باشد.





5

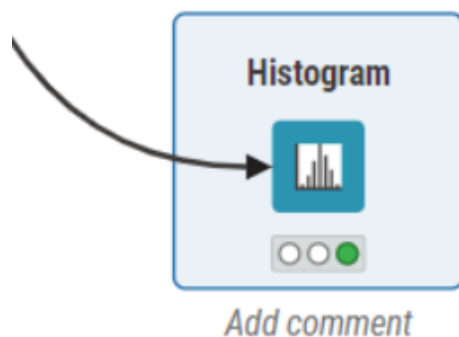
برای تعریف یک ویژگی جدید از Rule Engine node استفاده می‌کنیم . و در صورتی که دو ویژگی SibSP و Parch هر دو برابر با صفر بودند یعنی مسافر هیچ وابستگی به خانواده نداشته و تنهایی سفر کرده است.



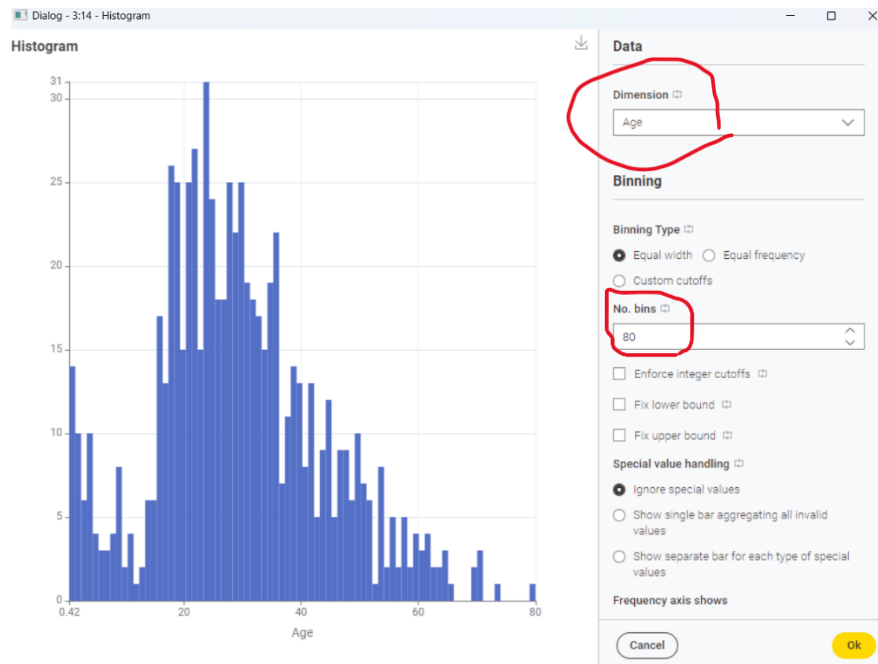
ass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	ALone
iber (inte...	String	String	Number (dou...	Number (inte...	Number (inte...	String	Number (dou...	String	String	Number (inte...
	Braund, Mr. Owl	male	22	1	0	A/5 21171	7.25	?	S	0
	Cumings, Mrs. .	female	38	1	0	PC 17599	71.283	C85	C	0
	Heikkinen, Miss	female	26	0	0	STON/O2. 3101	7.925	?	S	1
	Futrelle, Mrs. Ja	female	35	1	0	113803	53.1	C123	S	0
	Allen, Mr. Williar	male	35	0	0	373450	8.05	?	S	1
	Moran, Mr. Jam	male	?	0	0	330877	8.458	?	Q	1
	McCarthy, Mr. T	male	54	0	0	17463	51.862	E46	S	1
	Palsson, Mastei	male	2	3	1	349909	21.075	?	S	0
	Johnson, Mrs. C	female	27	0	2	347742	11.133	?	S	0
	Nasser, Mrs. Nik	female	14	1	0	237736	30.071	?	C	0

6

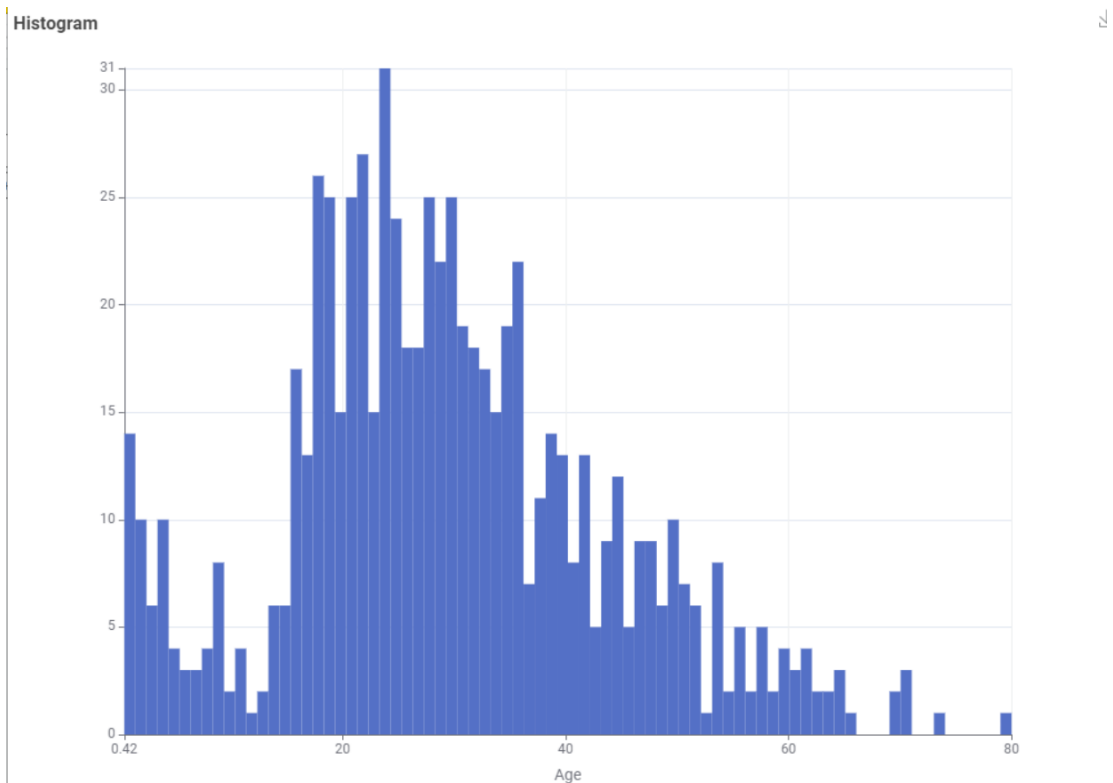
برای رسم histogram روی age باید از نمودار histogram استفاده کنیم که ورودی های خودش را از داده های csv تایید می‌گیرد.



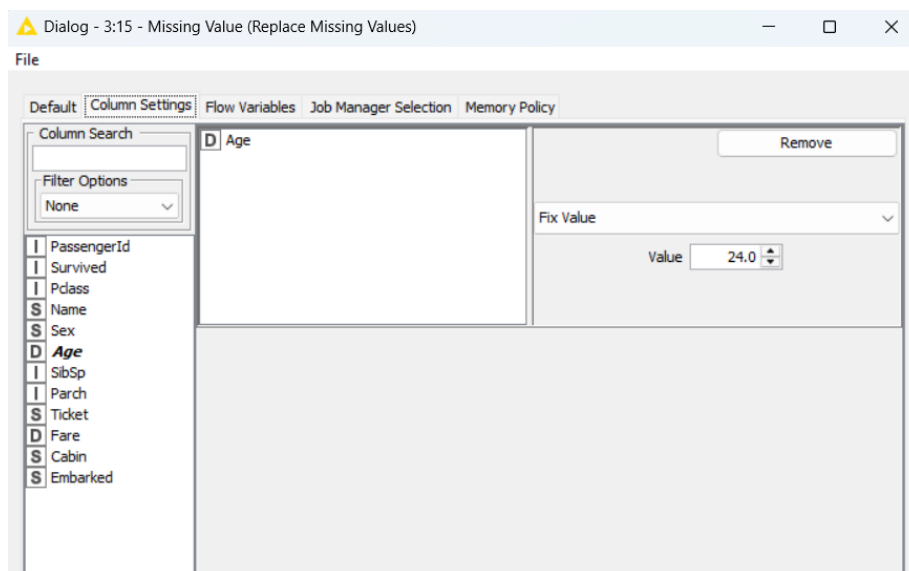
برای رسم هیستوگرام صرفاً موارد علامت زده شده در تصویر را تغییر داده ایم و باقی موارد مقدار default خود برنامه هستند



همانطور که در نمودار بدست آمده مشخص است، بیشتر افراد در بازه سنی بین 16 تا 36 سال هستند. تعداد افراد از 50 به بالا نسبت به باقی بازه سنی ها کمتر می باشد. بیشترین تعداد افراد با یک سن مشابه 31 نفر هست که نشان برابری با 24 سال است.



برای این سوال لازم است که از نود Missing Value استفاد کنیم. از سوال قبل و با استفاده از نمودار هیستوگرامی که روی Age بدست آوریم ، فهمیدیم که بیشترین سنی که در میان مسافران وجود داشته است برابر با 24 بوده پس مقدار 24 را جایگزین مقادیر Missing در ستون Age می‌کنیم.



مقادیر Age قبل از جایگزین کردن Missing value ها

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Braund, M	male	22	1	0	A/5 21171	7.25		S
2	1	1	Cumings, N	female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	Heikkinen, female		26	0	0	STON/O2.	7.925		S
4	1	1	Futrelle, M	female	35	1	0	113803	53.1	C123	S
5	0	3	Allen, Mr. J	male	35	0	0	373450	8.05		S
6	0	3	Moran, Mr	male		0	0	330877	8.4583		Q
7	0	1	McCarthy, male		54	0	0	17463	51.8625	E46	S
8	0	3	Palsson, M	male	2	3	1	349909	21.075		S
9	1	3	Johnson, N	female	27	0	2	347742	11.1333		S
10	1	2	Nasser, Mr	female	14	1	0	237736	30.0708		C
11	1	3	Sandstrom	female	4	1	1	PP 9549	16.7	G6	S
12	1	1	Bonnell, M	female	58	0	0	113783	26.55	C103	S
13	0	3	Saunders	male	20	0	0	A/5. 2151	8.05		S
14	0	3	Andersson	male	39	1	5	347082	31.275		S
15	0	3	Vestrom, N	female	14	0	0	350406	7.8542		S
16	1	2	Hewlett, N	female	55	0	0	248706	16		S
17	0	3	Rice, Mast	male	2	4	1	382652	29.125		Q
18	1	2	Williams, N	male		0	0	244373	13		S
19	0	3	Vander Pla	female	31	1	0	345763	18		S
20	1	3	Masselmani	female		0	0	2649	7.225		C

مقادیر Age بعد از جایگزین کردن Missing value ها

6	Row5	6	0	3	Moran, Mr. Jam	male	24
7	Row6	7	0	1	McCarthy, Mr. T	male	54
8	Row7	8	0	3	Palsson, Master	male	2
9	Row8	9	1	3	Johnson, Mrs. C	female	27
10	Row9	10	1	2	Nasser, Mrs. Nik	female	14
11	Row10	11	1	3	Sandstrom, Mis	female	4
12	Row11	12	1	1	Bonnell, Miss. E	female	58
13	Row12	13	0	3	Saunderscock, M	male	20
14	Row13	14	0	3	Andersson, Mr.	male	39
15	Row14	15	0	3	Vestrom, Miss.	female	14
16	Row15	16	1	2	Hewlett, Mrs. (M	female	55
17	Row16	17	0	3	Rice, Master. Eu	male	2
18	Row17	18	1	2	Williams, Mr. Ch	male	24
19	Row18	19	0	3	Vander Planke, I	female	31
20	Row19	20	1	3	Masselmani, Mi	female	24

8

برای پر کردن مقادیر missing در ویژگی کابین روش هایی که می توان استفاده کرد به این شکل است:

1.

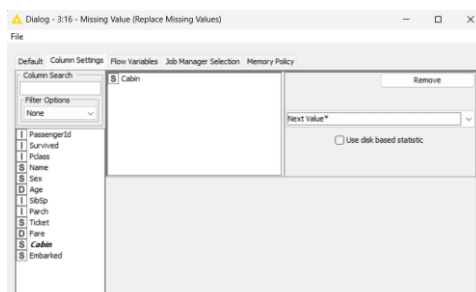
استفاده از روش های یادگیری ماشین: برای این کار می توان از ویژگی هایی مثل Pclass, Fare, Embarked استفاده کرد تا متوجه شویم مسافر در کدام کابین بوده است (مخصوصا چون حرف اول کابین براسا شماره عرشه کشتی می باشد می توان با این روش ، اینکه مسافر در کدام عرشه می باشد را حدس زد). به عنوان مثال اگر Embarked برای مسافری Q باشد و Fare آن از 20 کمتر باشد و Pclass آن 3 باشد احتمالا کابین آن بین حروف F و یا G بوده است.

2.

روش دیگر استفاده از نزدیک ترین همسایگی می باشد، باید بر اساس اطلاعات سایر مسافرین که اطلاعات آن ها نزدیک به مسافر فعلی مورد بررسی می باشد، تخمین بزنیم که شماره کابین مسافر چند بوده است.

3.

یک حالت این است که مقادیر را با مقدار های بعدی و یا قبلی خودشان جایگزین کنیم که این روش میتواند تا حدودی منطقی تر باشد اگر تعداد داده های گم شده خیلی زیاد نباشد که باعث شود اغلب فیلد های این ستون مقداری مشابه بگیرند.



Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
1	Cummings, Mrs. J.	female	38	1	0	PC 17599	71.283	C85
3	Heikkinen, Miss.	female	26	0	0	STON/O2. 3101	7.925	G6
1	Futrelle, Mrs. J.	female	35	1	0	113803	53.1	C123
3	Allen, Mr. William	male	35	0	0	373450	8.05	G6
3	Moran, Mr. James	male	?	0	0	330877	8.458	G6
1	McCarthy, Mr. Thomas	male	54	0	0	17463	51.862	E46
3	Palsson, Master G.	male	2	3	1	349909	21.075	G6
3	Johnson, Mrs. Charlotte	female	27	0	2	347742	11.133	G6
2	Nasser, Mrs. Nicholas	female	14	1	0	237736	30.071	G6
3	Sandstrom, Miss. Margaretha	female	4	1	1	PP 9549	16.7	G6
1	Bonnell, Miss. Elizabeth	female	58	0	0	113783	26.55	C103
3	Saunderscock, Mr. Thomas	male	20	0	0	A/5. 2151	8.05	G6
3	Andersson, Mr. Olof	male	39	1	5	347082	31.275	D56

4.

حالت بعدی استفاده از روش ها آماری همچون بیشترین مقداری که تکرار شده که البته این روش مناسب نیست به دلیل اینکه نتایج از حالت واقعی خارج می شود و نمی تواند چند نفر درون یک کابین باشند.

Dialog - 3:16 - Missing Value (Replace Missing Values)

File

Default Column Settings Flow Variables Job Manager Selection Memory Policy

Column Search

Filter Options

None

S Cabin

Remove

Most Frequent Value

PassengerId

Survived

Pclass

Name

Sex

Age

SibSp

Parch

Ticket

Fare

Cabin

Embarked

Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
Braund, Mr. Owen	male	22	1	0	A/5 21171	7.25	G6	S
Cummings, Mrs. J.	female	38	1	0	PC 17599	71.283	C85	C
Heikkinen, Miss.	female	26	0	0	STON/O2. 3101	7.925	G6	S
Futrelle, Mrs. J.	female	35	1	0	113803	53.1	C123	S
Allen, Mr. William	male	35	0	0	373450	8.05	G6	S
Moran, Mr. James	male	?	0	0	330877	8.458	G6	Q
McCarthy, Mr. Thomas	male	54	0	0	17463	51.862	E46	S
Palsson, Master G.	male	2	3	1	349909	21.075	G6	S
Johnson, Mrs. Charlotte	female	27	0	2	347742	11.133	G6	S
Nasser, Mrs. Nicholas	female	14	1	0	237736	30.071	G6	C
Sandstrom, Miss. Margaretha	female	4	1	1	PP 9549	16.7	G6	S
Bonnell, Miss. Elizabeth	female	58	0	0	113783	26.55	C103	S
Saunderscock, Mr. Thomas	male	20	0	0	A/5. 2151	8.05	G6	S
Andersson, Mr. Olof	male	39	1	5	347082	31.275	G6	S
Vestrom, Miss. Anna	female	14	0	0	350406	7.854	G6	S
Hewlett, Mrs. (Mrs. Mary)	female	55	0	0	248706	16	G6	S
Rice, Master Eugene	male	2	4	1	382652	29.125	G6	Q

2.1 نقد محدودیت های روش های بصری سازی

1

خیر روش های سنتی مثل هیستوگرام و Scatter plot روش های مناسبی نیستند، چرا که تعداد داده و ویژگی برای این مثلا بسیار زیاد می باشد، مثلا از آنجایی که روش Scatter plot می تواند نهایتا رابطه بین دو الی سه attribute را نشان دهد، پس برای مقایسه 50 ویژگی نمی توان از این روش استفاده کرد، چرا که مقایسه هر 50 داده باهم ممکن نیست و ثانياً برای مقایسه دودویی هر روش وقت زیاد با دقت پایین حاصل می شود. از طرفی به دلیل داشتن 1000 نمونه روش هیستوگرام هم روش مناسبی نیست مگر اینکه برای تک تک نمونه های داخل attribute یک bine در نظر بگیریم و صرفاً یک بازه را مدنظر قرار دهیم، و اگر بخواهیم هیستوگرام را برای یک attribute بدست آوریم، باید 50 تا نمودار بکشیم که مقدار زیادی می باشد.

2

اگر بخواهیم که در فضای دوبعدی نشان دهیم می توانیم از تکنیک های کاهش ابعاد مثل (Principal Component Analysis) استفاده کنیم که چند مولفه را بدیل به دو یا سه مولفه می کند و در این حالت می توانیم نمودار Scatter plot ان را رسم کنیم.

همچنین برای بررسی همبستگی میان ویژگی ها می توان از نمودار Heatmap هم استفاده کرد. همچنین برای چنین داده هایی می توان از نمودار های جدیدتری همچون Parallel Coordinates استفاده کرد، که هر متغیر را می توان با یک شدت رنگ تعریف کرد.

3 تحلیل فاصله ها در مجموعه داده

3.1

1

ابتدا باید جدول گفته شده را توسط نود Table Creator بسازیم.

The screenshot shows the Orange3 data mining software interface. In the center, there is a 'Table Creator' widget. To its right, a 'Table Creator Settings' dialog box is open, showing a table with 10 rows and 3 columns (x1, x2, x3). The main window below the dialog shows the resulting table with 4 rows and 3 columns (x1, x2, x3).

#	RowID	x1 Number (integer)	x2 Number (integer)	x3 Number (integer)
1	point1	2	3	4
2	point2	1	5	6
3	point3	4	6	7
4	point4	3	2	8

Distance Euclidean .1

برای محاسبه فاصله Euclidean نیاز داریم تا از نود numeric distance استفاده کنیم و از نود distance matrix calculate هم برای تشکیل ماتریس فاصله استفاده می‌کنیم.

File

Distance Configuration | Flow Variables | Job Manager Selection

☒ Manual Selection ☐ Wildcard/Regex Selection

Exclude

Filter

No columns in this list

☐ Enforce exclusion

Include

Filter

x1
x2
x3

☒ Enforce inclusion

Distance Selection

Standard Distance (Euclidean/ Manhattan)

Configuration

☒ Euclidean

☐ Manhattan

☐ Maximum

☐ Custom 'p' 2.0

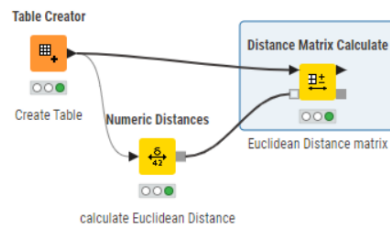
☐ Normalize distance (Requires normalized input vectors)

Missing Values

☒ Fail (fails if a missing value cell occurs during computation.)

☐ Assume equal (Assume a missing value has the value of the respective counterpart - this will add 0 to the sum of pair wise absolute differences)

☐ Average distance (i.e. ignores the missing value and the corresponding value.)



► 1: Table containing distance matrix column ■ 2: Matrix Distance Measure ⚙ Flow Variables

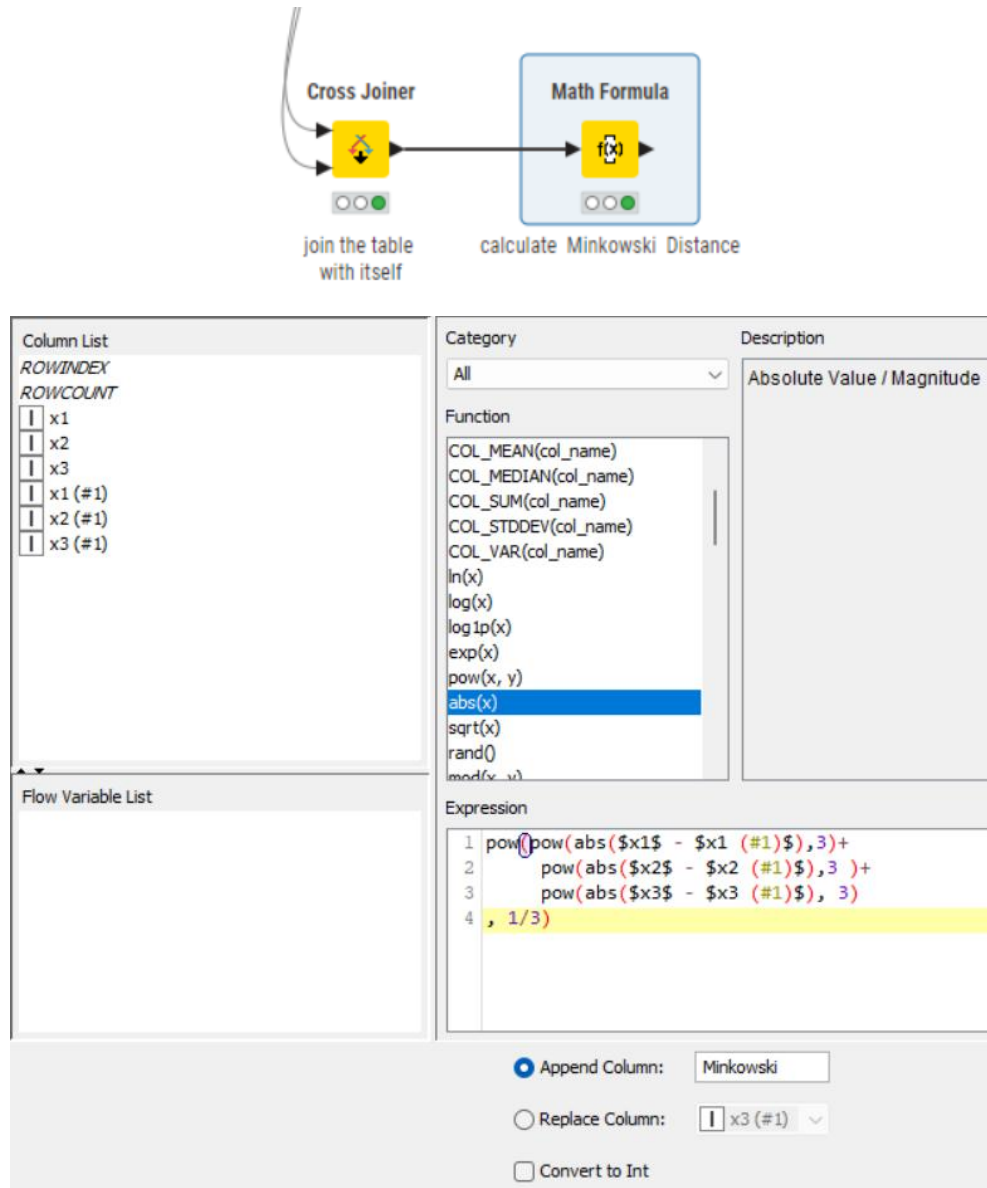
Rows: 4 | Columns: 4

Table | Statistics

#	RowID	x1 Number (integer)	x2 Number (integer)	x3 Number (integer)	Distance Distance vector
1	point1	2	3	4	0 []
2	point2	1	5	6	1 [3.0]
3	point3	4	6	7	2 [4.69041575982343, 3.31662479035]
4	point4	3	2	8	3 [4.242640687119285, 4.1231056256]

2. Distance Minkowski با $r=3$

چون به صورت مستقیم نودی وجود ندارد که بتوانیم این فاصله را بدست آوریم، باید به صورت دستی خودمان این روش را پیاده سازی کنیم. پس برای این کار لازم است که ابتدا از نود `crosse join` استفاده کنیم تا `table` ساخته شده را با خودش `join` کنیم تا بتوانیم با توجه به `table` بدست آمده، فاصله هر نقطه با نقاط دیگر را بدست آوریم. پس برای اینکار از فرمول رابطه Minkowski استفاده می‌کنیم.

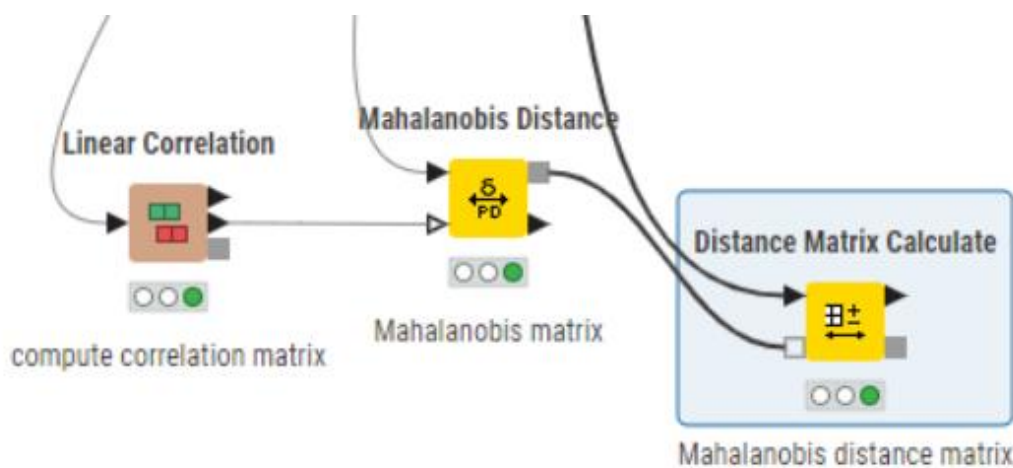


با توجه به شکل زیر هر کدام از سطر های ستون Minkowski به ترتیب می‌شود: فاصله بین نقطه اول با خودش، فاصله نقطه اول با نقطه دوم، فاصله نقطه اول با نقطه سوم و فاصله نقطه اول با نقطه چهارم و بعد از آن فاصله نقطه دوم تا خودش و

#	RowID	x1 Number (integer)	x2 Number (integer)	x3 Number (integer)	x1 (#1) Number (integer)	x2 (#1) Number (integer)	x3 (#1) Number (integer)	Minkowski Number (double)
1	point1_	2	3	4	2	3	4	0
2	point1_	2	3	4	1	5	6	2.571
3	point1_	2	3	4	4	6	7	3.958
4	point1_	2	3	4	3	2	8	4.041
5	point2_	1	5	6	2	3	4	2.571
6	point2_	1	5	6	1	5	6	0
7	point2_	1	5	6	4	6	7	3.072
8	point2_	1	5	6	3	2	8	3.503
9	point3_	4	6	7	2	3	4	3.958
10	point3_	4	6	7	1	5	6	3.072
11	point3_	4	6	7	4	6	7	0
12	point3_	4	6	7	3	2	8	4.041
13	point4_	3	2	8	2	3	4	4.041
14	point4_	3	2	8	1	5	6	3.503
15	point4_	3	2	8	4	6	7	4.041
16	point4_	3	2	8	3	2	8	0

Distance Mahalanobis.2

در ابتدا باید correlation را حساب کنیم و بعد با استفاده از node های Mahalanobis و Distance matrix calculate فاصله بین نقاط را بدست آوریم.



نهایتا خروجی به شکل زیر می باشد:

#	RowID	x1 Number (integer)	x2 Number (integer)	x3 Number (integer)	Distance Distance vector
1	point1	2	3	4	0 []
2	point2	1	5	6	1 [3.978684573661562]
3	point3	4	6	7	2 [4.2426637922478925, 3.122781860]
4	point4	3	2	8	3 [4.284686258373692, 3.9579678054]

Euclidean Distance

از نقطه	به نقطه 1	به نقطه 2	به نقطه 3	به نقطه 4
نقطه 1	0	3	4.69041575982343	4.242640687119285
نقطه 2	3	0	3.3166247903554	
نقطه 3	4.69041575982343	3.3166247903554	0	
نقطه 4	4.2426406871192	4.12310562561766	4.242640687119285	0

Minkowski Distance

از نقطه	به نقطه 1	به نقطه 2	به نقطه 3	به نقطه 4
نقطه 1	0	2.571	3.958	4.041
نقطه 2	2.571	0	3.072	3.503
نقطه 3	3.958	3.072	0	4.041
نقطه 4	4.041	3.503	4.041	0

Mahalanobis Distance

از نقطه	به نقطه 1	به نقطه 2	به نقطه 3	به نقطه 4
نقطه 1	0	3.97868457366156 2	4.24266379224789 25	4.28468625837369 2
نقطه 2	3.97868457366156 2	0	3.12278186072464	3.95796780545605 34
نقطه 3	4.24266379224789 25	3.12278186072464	0	4.28089851927698 7
نقطه 4	4.28468625837369 2	3.95796780545605 34	4.28089851927698 7	0

اگر همبستگی داده ها بهم کم باشد میتوان از دو ویژگی Minkowski و Euclidean استفاده کرد، اما در اصل Minkowski میتواند فاصله های بیشتر بین داده ها را به خوبی نشان بدهد، پس میتواند معیار

بهتری برای نشان دادن تفاوت های میان داده ها باشد، اما اگر داده ها با توجه به شرایط گفته شده در بالا همبستگی داشته باشند بایکدیگر این روش بهتر است Mahalanobis.

3.3 با توجه به چهار زمینه داده ای بالا:

1

در این حالت بهتر است که از فاصله Minkowski با $r=1$ استفاده کنیم، زیرا که اگر داده ها همبسته نباشند ممکن است که Mahalanobis نتایج غلط باشد (اگر همبسته باشند گزینه خوبی هست) و فاصله اقلیدسی هم برای داده های گسسته جواب بی معنی ممکن است بدهد.

2

فاصله Mahalanobis بهترین گزینه است چون رابطه همبستگی ویژگی ها را در نظر می گیرد و داده ها را در فضایی نرمال شده مقایسه می کند.

3

فاصله اقلیدسی و Minkowski نیاز به نرمال سازی داده ها دارد. اما Mahalanobis به دلیل در نظر گرفتن کواریانس، نیاز به نرمال سازی ندارد و داده ها را متعادل می کند.

4

اگر داده ها به شدت غیر نرمال باشند فاصله Minkowski خیلی بهتر از دو فاصله دیگر نتایج را برمی گرداند.