

Intelligent Modeling Introduction

YModel



www.raqsoft.com.cn

CONTENTS

01

Data source

02

Data exploration

03

Preprocessing

04

Modeling

05

**Model
performance**

06

Prediction

07

**Integration
solution**

CONTENTS

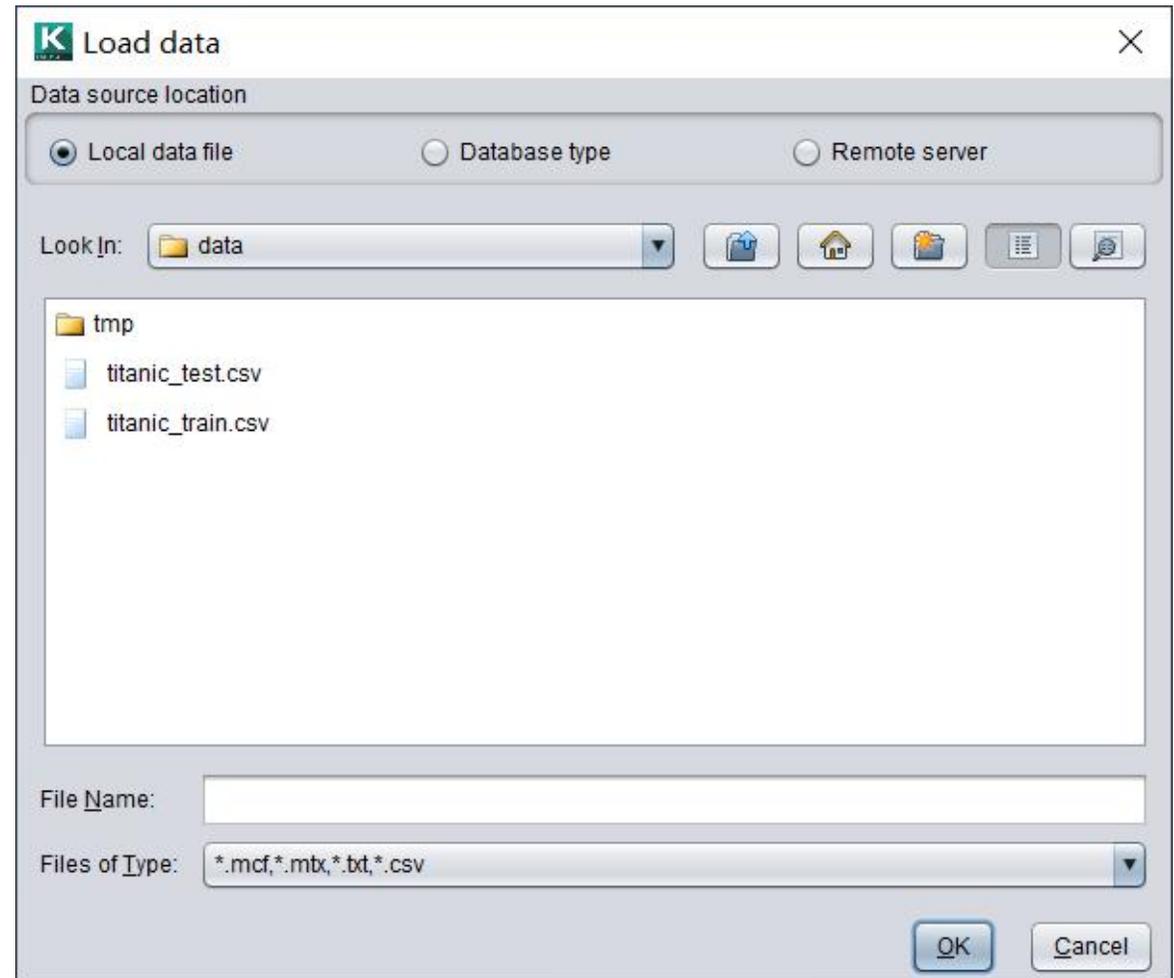
1. Local data file
2. Database

01

Data Source

➤ 1. Local data file

Intelligent modeling supports TXT, CSV and other data files.





1. Local data file



After selecting a file, you can define the parameter configuration of the data file.

The screenshot shows the 'Load data' dialog box with the following configuration options:

- Create data file name: `titanic_train.mtx`
- Import the first line as variable name
- Omit all quotation marks
- Check Column Count
- Delete a line when column count does not match value count at line 1
- Use double quotation marks as escape characters
- Delimiter: `,`
- Charset: `GBK`
- Date format: `yyyy/MM/dd`
- Time format: `HH:mm:ss`
- Date time format: `yyyy/MM/dd HH:mm:ss`
- Locale: `English`
- Missing values (bar-separated): `NULL|N/A`

The 'Preview data' section shows a table with the following columns: PassengerId, Survived, Pclass, Name, Sex, and Age. The table displays 21 rows of data:

PassengerId	Survived	Pclass	Name	Sex	Age
1	0	3	Braund, Mr. Owen Harris	male	
2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	
3	1	3	Heikkinen, Miss. Laina	female	
4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	
5	0	3	Allen, Mr. William Henry	male	
6	0	3	Moran, Mr. James	male	
7	0	1	McCarthy, Mr. Timothy J	male	
8	0	3	Paisson, Master. Gosta Leonard	male	
9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	
10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	
11	1	3	Sandstrom, Miss. Marguerite Rut	female	
12	1	1	Bonnell, Miss. Elizabeth	female	
13	0	3	Saunderscock, Mr. William Henry	male	
14	0	3	Andersson, Mr. Anders Johan	male	
15	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	
16	1	2	Hewlett, Mrs. (Mary D Kingcome)	female	
17	0	3	Rice, Master. Eugene	male	
18	1	2	Williams, Mr. Charles Eugene	male	
19	0	3	Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele)	female	
20	1	3	Masselmani, Mrs. Fatima	female	
21	0	2	Fynney, Mr. Joseph J	male	

At the bottom of the dialog, there are buttons for 'Cancel', '< Previous', 'Next >', and 'Finish'.



➤ 1. Local data file

Next, you can define the variable type, date format, and selection status.

Variable types can be automatically detected or be configured by importing the data dictionary.

The format of data dictionary is as follows:

Name	Type	DateFormat	Used	Importance
PassengerId	Identity		TRUE	0
Survived	Binary		TRUE	0
Pclass	Categorical		TRUE	0
Name	Text		FALSE	0
Sex	Binary		TRUE	0
Age	Numerical		TRUE	0
SibSp	Categorical		TRUE	0
...

K Load data

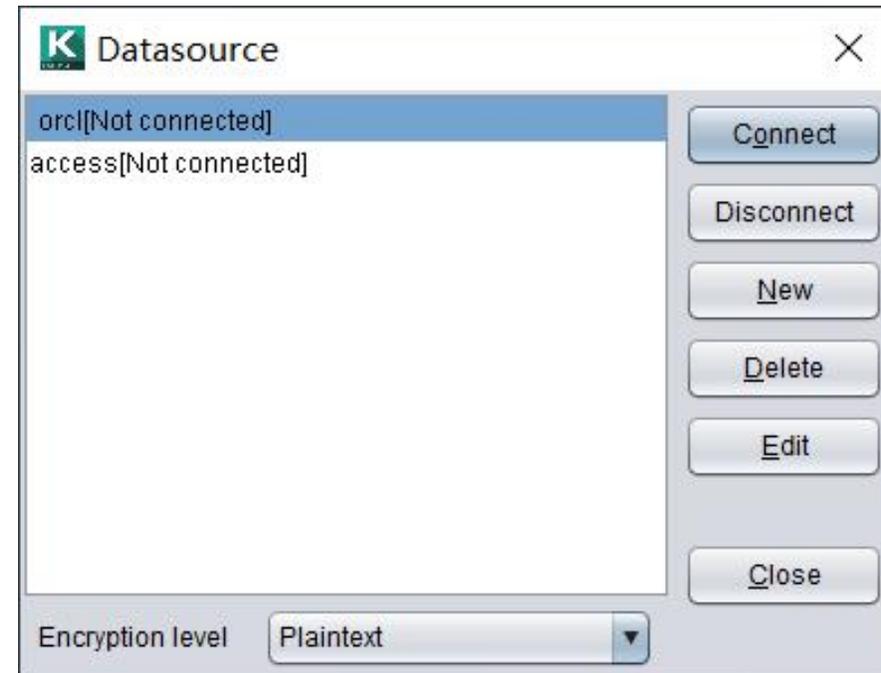
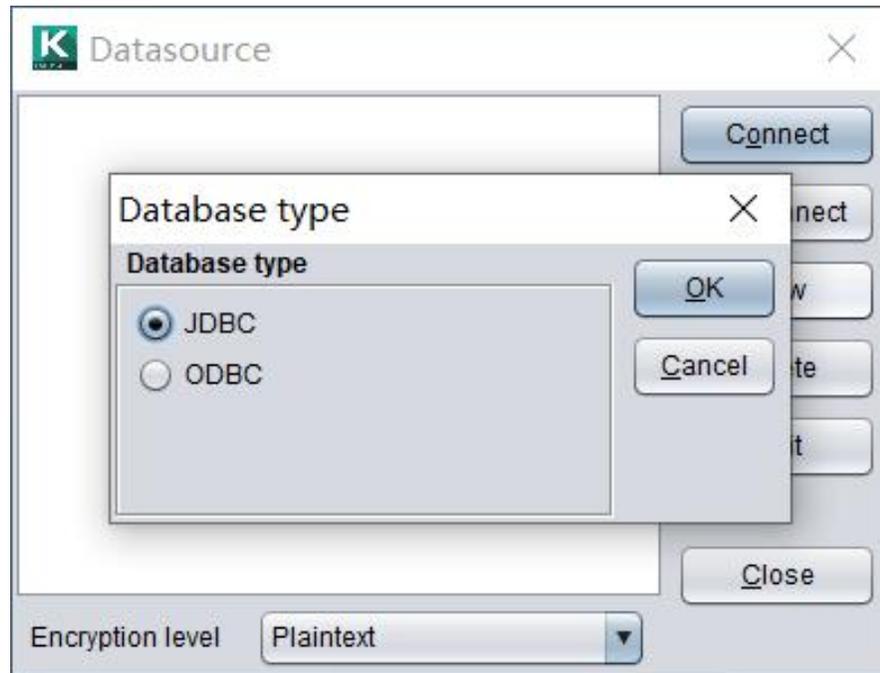
Import data dictionary Note: Unselected variables won't be imported.

NO.	Variable name	Type	Date format	<input checked="" type="checkbox"/> Select
1	PassengerId	Automatic		<input checked="" type="checkbox"/>
2	Survived	Automatic		<input checked="" type="checkbox"/>
3	Pclass	Automatic		<input checked="" type="checkbox"/>
4	Name	Automatic		<input checked="" type="checkbox"/>
5	Sex	Automatic		<input checked="" type="checkbox"/>
6	Age	Automatic		<input checked="" type="checkbox"/>
7	SibSp	Automatic		<input checked="" type="checkbox"/>
8	Parch	Automatic		<input checked="" type="checkbox"/>
9	Ticket	Automatic		<input checked="" type="checkbox"/>
10	Fare	Automatic		<input checked="" type="checkbox"/>
11	Cabin	Automatic		<input checked="" type="checkbox"/>
12	Embarked	Automatic		<input checked="" type="checkbox"/>

2. Database



In the data source window, you can define two data source connections: JDBC and ODBC.



➤ 2. Database



JDBC Datasource

The screenshot shows the 'Datasource' dialog box with the 'General properties' tab selected. The fields are filled with the following information:

- Datasource name: orcl
- Database vendor: ORACLE
- Driver: oracle.jdbc.driver.OracleDriver
- Datasource URL: jdbc:oracle:thin:@127.0.0.1:1521:orcl
- User: System
- Password: ****
- Batch size: 0
- Options: Qualify object with schema, Enclose object name in quotes

ODBC Datasource

The screenshot shows the 'ODBC datasource' dialog box. The fields are filled with the following information:

- Datasource name: access
- ODBC name: (empty)
- Username: (empty)
- Password: (empty)
- Options: Qualify object with schema, Case sensitive, Enclose object name in quotes

➤ 2. Database



Next, you can use the configured data source to edit the SQL statement for data loading.

Load data [X]

Data source location

Local data file Database type Remote server

Create data file name: scores.mtx

Table | Field | Where | Group | Having | Sort | SQL

Available table

SCORES

Selected table

SCORES

Data source: orcl Schema: WN

OK Cancel

Load data [X]

Data source location

Local data file Database type Remote server

Create data file name: scores.mtx

Table | Field | Where | Group | Having | Sort | SQL

SELECT * FROM SCORES

Data source: orcl Schema: WN

OK Cancel

CONTENTS

1. Basic characteristics
2. Statistics of discrete variables
3. Continuous variable statistics
4. Data exploration report
5. Data quality report

02

Data Exploration



➤ 1. Basic characteristics

After importing the data, the basic characteristics of the data are displayed:

The target variable is survived (it needs to be set by the user), with 12 variables and 891 records.

Automatically parses the types of each variable and the recommended selection status.

Model file: titanic_train.pcf | Model performance | Model presentation | Model options

Data file: titanic_train.mtx | Reload data

Target variable: **Survived** | Set | Variable filter | ↑ | ↓

NO.	Variable name	Type	Date format	Select
1	PassengerId	ID		<input checked="" type="checkbox"/>
2	Survived	Binary variable		<input checked="" type="checkbox"/>
3	Pclass	Categorical variable		<input checked="" type="checkbox"/>
4	Name	ID		<input type="checkbox"/>
5	Sex	Binary variable		<input checked="" type="checkbox"/>
6	Age	Numerical variable		<input checked="" type="checkbox"/>
7	SibSp	Categorical variable		<input checked="" type="checkbox"/>
8	Parch	Categorical variable		<input checked="" type="checkbox"/>
9	Ticket	Categorical variable		<input checked="" type="checkbox"/>
10	Fare	Numerical variable		<input checked="" type="checkbox"/>
11	Cabin	Categorical variable		<input checked="" type="checkbox"/>
12	Embarked	Categorical variable		<input checked="" type="checkbox"/>

Search variable: | Import 891 rows, 12 variables

➤ 1. Basic characteristics



The variable types of intelligent modeling are as follows:

Variable type	Description
Numerical variable	Variable with real number value
Single value variable	Variables containing only one category (excluding missing values)
Binary variable	Variables with only two categories (excluding missing values)
Count variable	Variable with natural value
Categorical variable	Variables with more than two classifications (excluding missing values)
ID	Unique identifier
Time and date	Date, time or datetime variable
Long text	Variables with a length of more than 128 bytes and a large number of classifications

The target variables of intelligent modeling support **binary variables, numerical variables, count variables and categorical variables.**

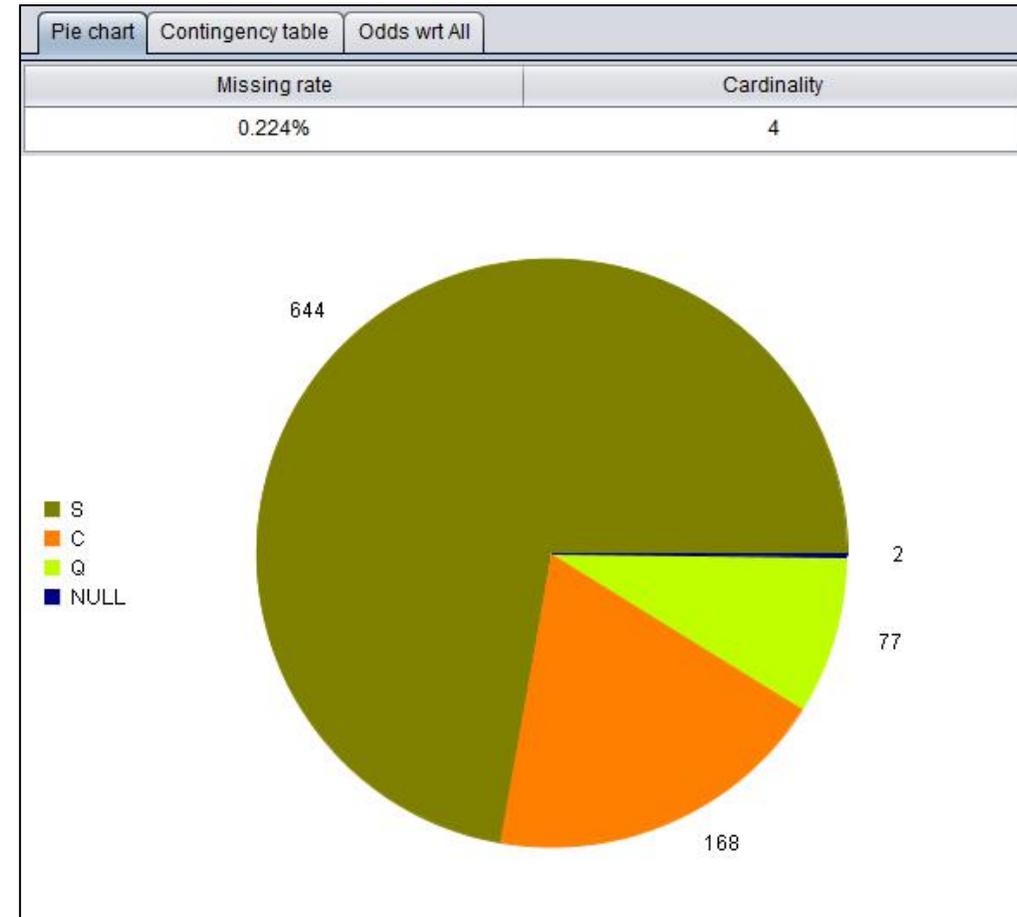
➤ 2. Statistics of discrete variables

Discrete variables include single value variables, binary variables and categorical variables.

Missing rate: the percentage of missing values in all data.

Potential: the number of members of the set that can be valued by a discrete variable.

Pie chart shows the proportion of each classification.



➤ 2. Statistics of discrete variables

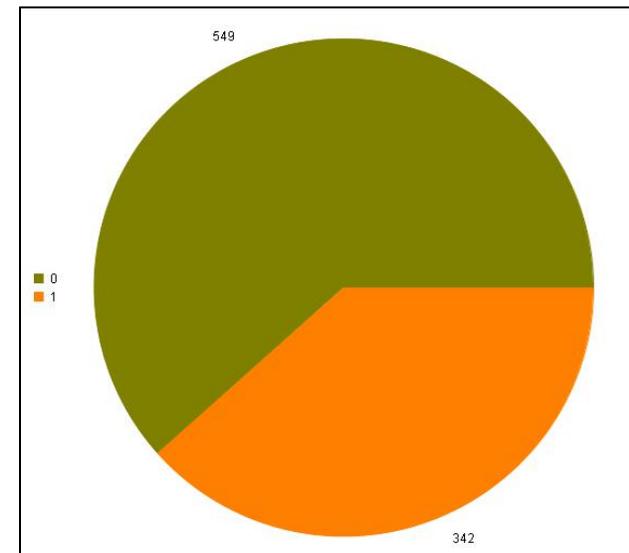


Target variable is binary variable: frequency distribution table of grouped target

In the frequency distribution table of grouped target, samples are grouped according to the classification value, and the number of samples in each group, the number of positive samples, the rate of positive samples and odds(occurrence ratio) are observed.

The positive sample of binary target variable refers to the classification value with a small number of samples. As can be seen from the right figure, in this example, the positive sample is a record with a target variable value of 1.

Categorical Level	Frequency	Positive Frequency	Positive Ratio	Odds wrt All
S	644	217	33.696%	0.878
C	168	93	55.357%	1.442
Q	77	30	38.961%	1.015
NULL	2	2	100%	2.605
All	891	342	38.384%	1



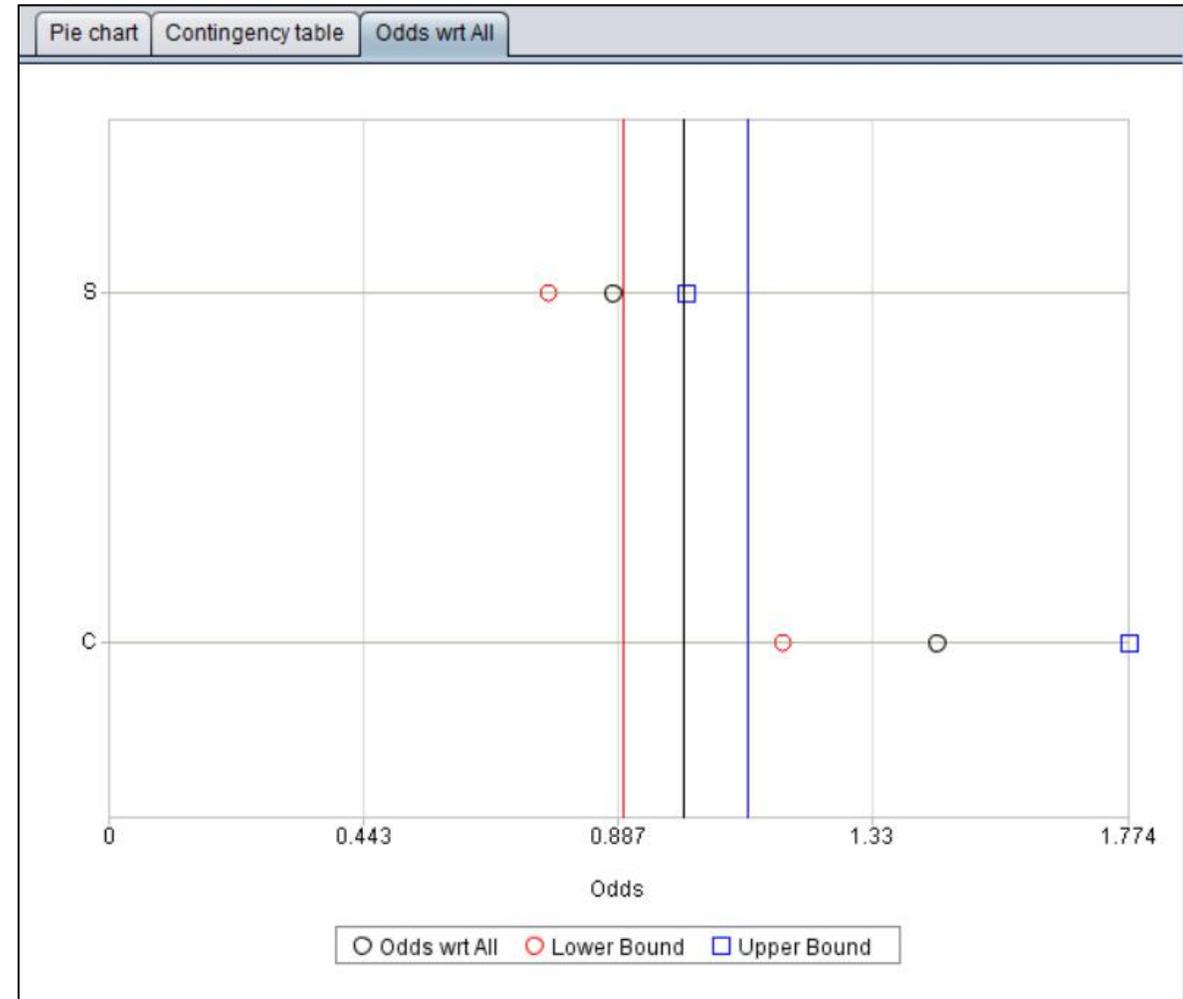
Pie chart of target variable

➤ 2. Statistics of discrete variables



Target variable is binary variable: frequency distribution table of grouped target

The odds wrt all graph shows the odds for each group of samples and the total odds. Classification with fewer samples (less than 100 samples) is not drawn.



➤ 2. Statistics of discrete variables



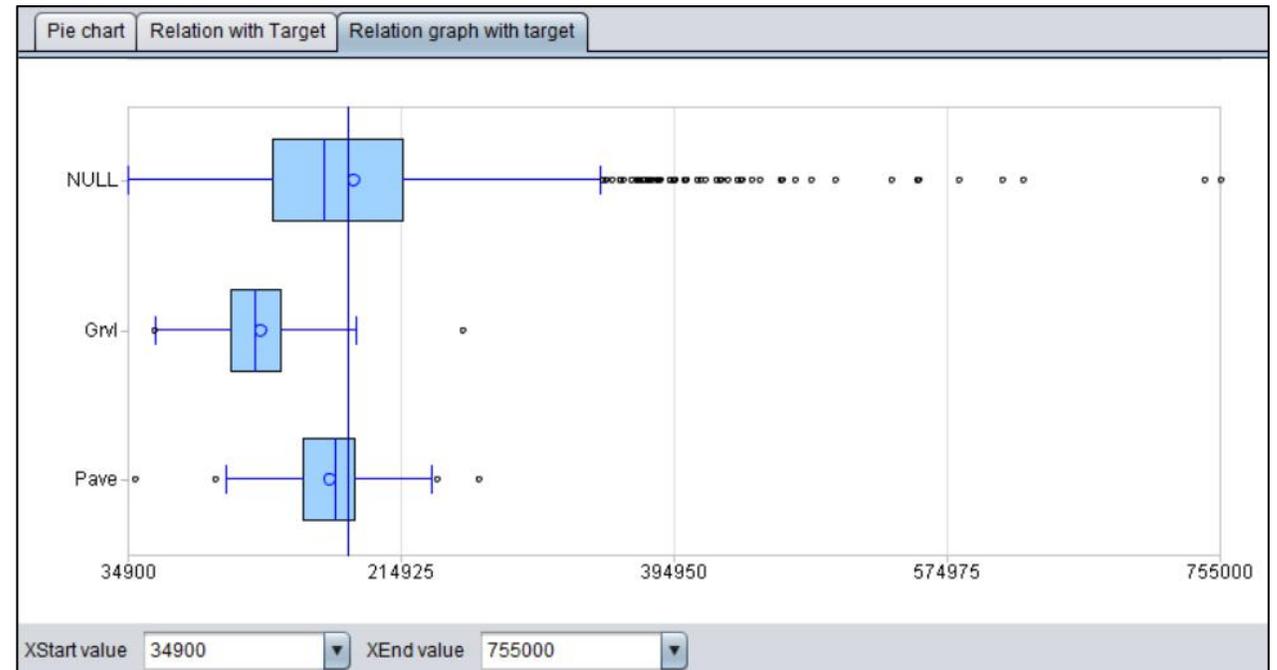
Target variable is numerical variable: statistics of grouped target, statistics of grouped target graph

Grouped target statistics group the samples according to the categorical value, and observe the statistics of each group of samples.

Including: frequency, average, standard deviation, median, minimum, maximum and Z-STAT.

The statistical graph of grouped target, in the form of box line chart, more intuitively represents the distribution of each group of samples. A box line chart can be used to mark outliers.

Categorical variable	Frequency	Average	Standard deviation	Median	Minimum	Maximum	Z-STAT
NULL	1369	183452.131	80667.145	165000	34900	755000	1.19
Grv	50	122219.08	34780.781	119000	52500	256000	-5.276
Pave	41	168000.585	38370.375	171900	40000	265979	-1.051



3. Continuous variable statistics



Continuous variables include numerical variables, count variables and time date variables.

Descriptive statistics show the basic statistical information of the data.

Frequency distribution diagram includes frequency distribution histogram, normal distribution curve and box line chart.

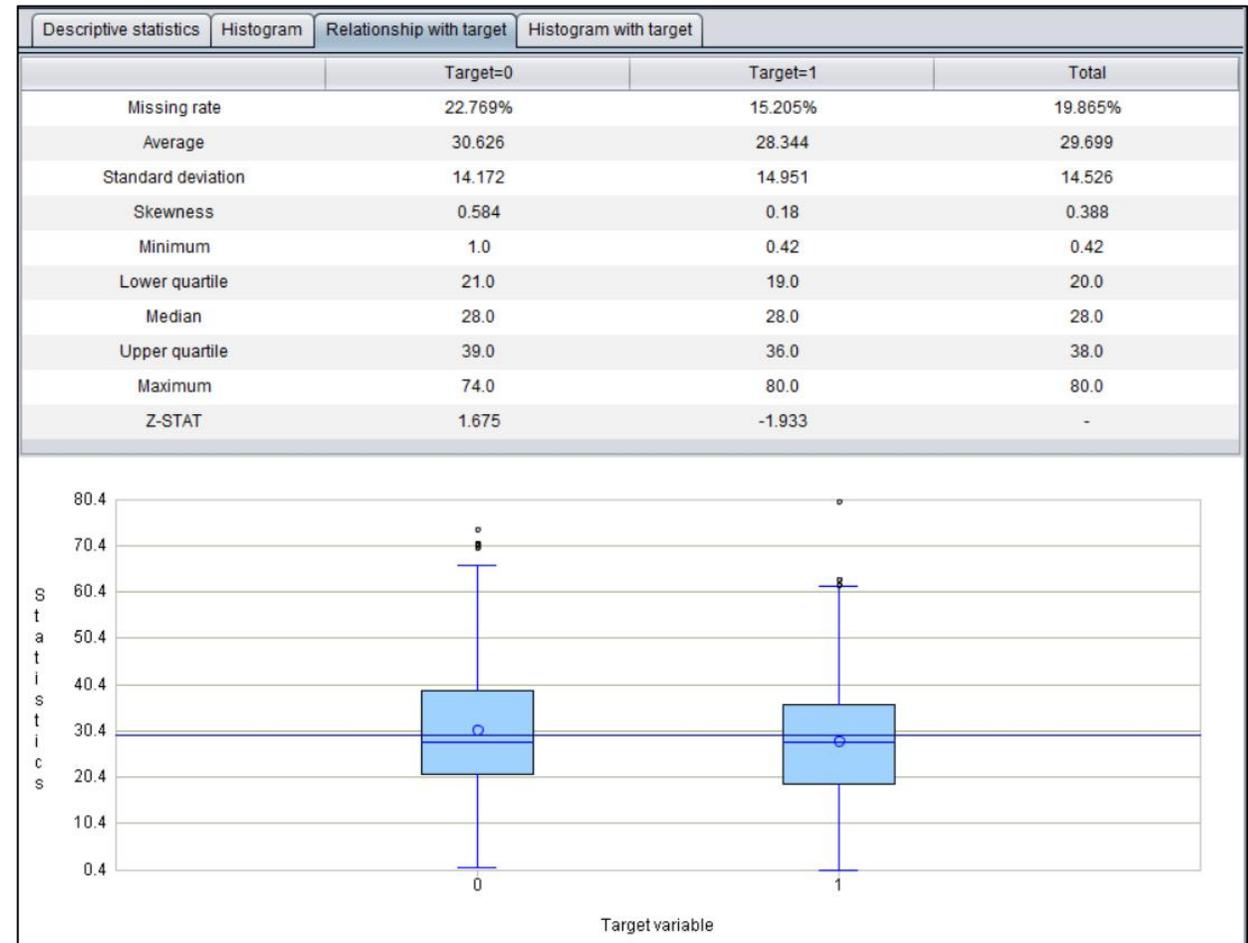


➤ 3. Continuous variable statistics



Target variable is binary variable: descriptive statistics of grouped target

Descriptive statistics of grouped target group the samples according to the target variable values, make statistics respectively, and draw the corresponding box line chart.

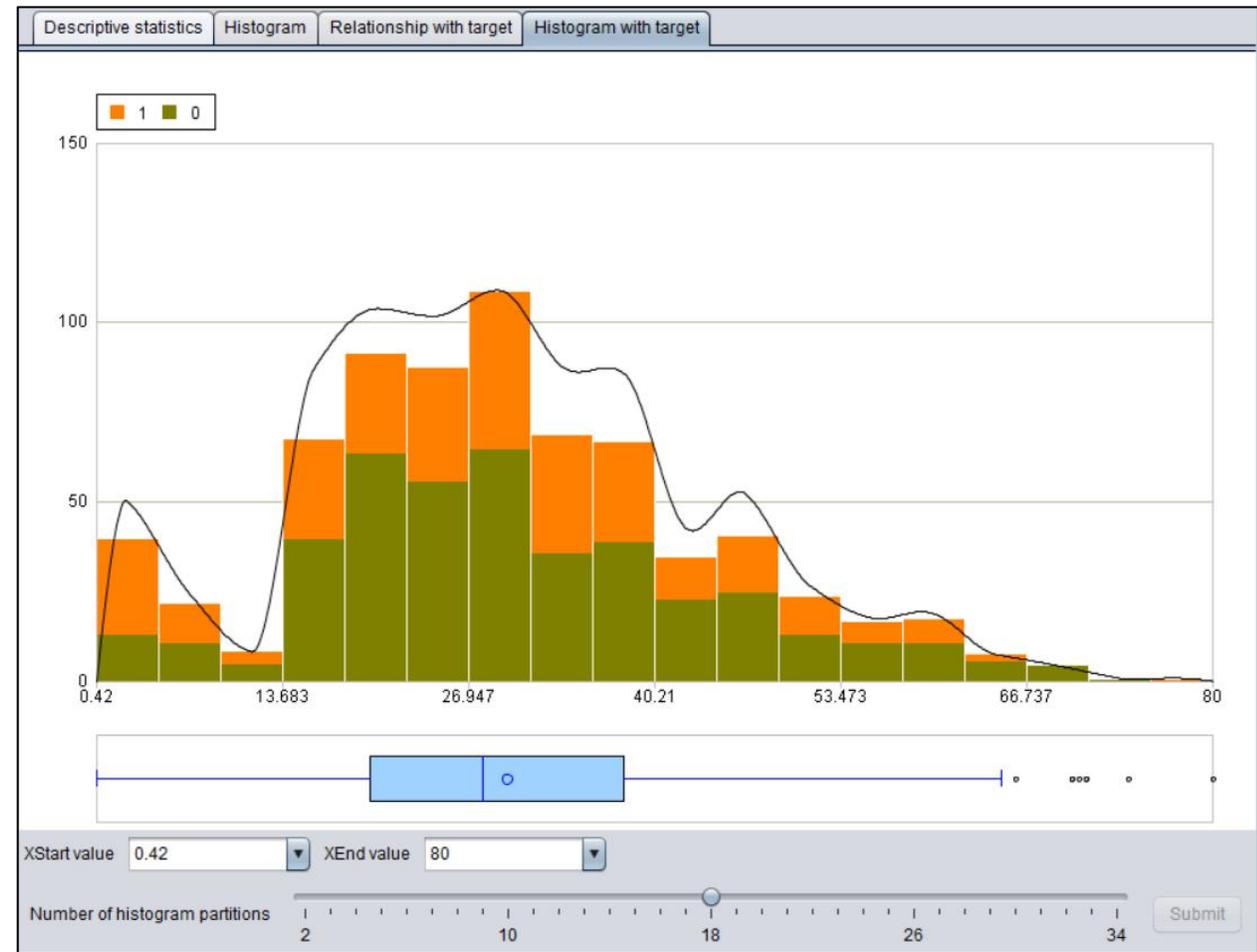


➤ 3. Continuous variable statistics



Target variable is binary variable : frequency distributions of grouped target

Frequency distributions of grouped target: the samples in each interval are grouped according to the target variable value, and the frequency is displayed in different colors.





➤ 3. Continuous variable statistics

Target variable is a numerical variable: target variable correlation coefficient

Pearson correlation coefficient: used to describe the linear correlation between two continuous variables.

Spearman rank correlation coefficient: used to describe the rank correlation between two continuous variables.

The greater the absolute value of the correlation coefficient, the greater the correlation between the two variables.

Descriptive statistics		Histogram		Correlation		Scatter Plot	
Pearson				Spearman			
0.7086				0.7313			

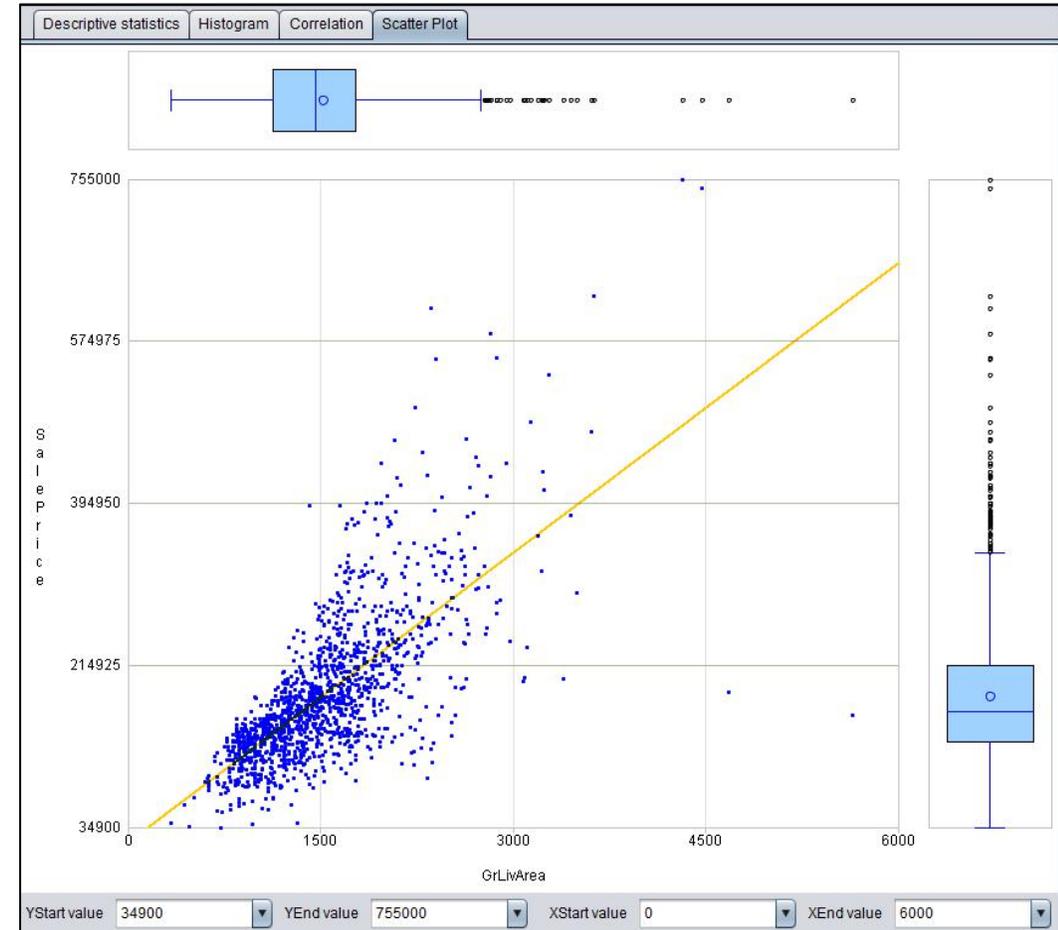
Above is the correlation coefficient between basement area and house price. It can be seen that there is a strong correlation between the two.

➤ 3. Continuous variable statistics



Target variable is a numerical variable : single factor scatter plot

The single factor scatter plot intuitively shows the correlation distribution of current variable (basement area) and target variable (house price). The yellow line is the regression line.





4. Data exploration report



Provide the function of exporting data exploration report to Excel file. Sheet1 is the basic information of variables:

	A	B	C	D	E	F	G	H	I	J
1	1.Numerical Variable:									
2		Variable Name	Minimum	Lower Quartile	Median	Mean	Upper Quartile	Maximum	Missing Rate	Skewness
3		Age	0.42	20.0	28.0	29.699	38.0	80.0	19.865%	0.388
4		Fare	0.0	7.896	14.454	32.204	31.0	512.329	0.0%	4.779
5	2.Count Variable: None									
6	3.Categorical Variable (Binary and Unary):									
7		Variable Name	Cardinality	Missing Rate						
8		Survived	2	0.0%						
9		Pclass	3	0.0%						
10		Sex	2	0.0%						
11		SibSp	7	0.0%						
12		Parch	7	0.0%						
13		Ticket	681	0.0%						
14		Cabin	148	77.104%						
15		Embarked	4	0.224%						
16	4.ID:									
17		Variable Name								
18		PassengerId								
19		Name								
20	5.Date Time: None									
21	6.Text String: None									
22										
23										
24										
25										
26										
27										
28										
29										
30										

Variables | With Binary Target Variable



4. Data exploration report



Sheet2 is the correlation between various variables and target variable.

	A	B	C	D	E	F	G	H	I	J	K	
1	1.Numerical Variable:											
2		Age										
3		Target	Frequency	Minimum	Lower Quartile	Median	Mean	Upper Quartile	Maximum	Missing Rate	Z_STAT wrt Overall	
4		1	342	0.42	19.0	28.0	28.344	36.0	80.0	15.205%	-1.933	
5		0	549	1.0	21.0	28.0	30.626	39.0	74.0	22.769%	1.675	
6												
7		Fare										
8		Target	Frequency	Minimum	Lower Quartile	Median	Mean	Upper Quartile	Maximum	Missing Rate	Z_STAT wrt Overall	
9		1	342	0.0	12.475	26.0	48.395	57.0	512.329	0.0%	6.232	
10		0	549	0.0	7.854	10.5	22.118	26.0	263.0	0.0%	-4.919	
11	2.Count Variable: None											
12	3.Categorical Variable (Binary, Unary):											
13		Pclass										
14		Category	Frequency	Positive Frequency	Positive Ratio	Odds wrt All						
15		1	216	136	62.963%	1.64						
16		2	184	87	47.283%	1.232						
17		3	491	119	24.236%	0.631						
18		All	891	342	38.384%	1.0						
19												
20		Sex										
21		Category	Frequency	Positive Frequency	Positive Ratio	Odds wrt All						
22		female	314	233	74.204%	1.933						
23		male	577	109	18.891%	0.492						
24		All	891	342	38.384%	1.0						
25												
26		SibSp										
27		Category	Frequency	Positive Frequency	Positive Ratio	Odds wrt All						
28		0	608	210	34.539%	0.9						
29		1	209	112	53.589%	1.396						
30		2	28	13	46.429%	1.21						

Variables

With Binary Target Variable



5. Data quality report



Provide the function of exporting data quality report to PDF file. Some contents are as follows:

These observations from 891 unique ID. Since the number of unique Id is equal to the number of observations, time sensitive information can not be studied using this data set.

Variables that have all “empty“ value are not exist.

Variables that have more than 99% of missing values are not exist.

Variables that have missing values between 95% and 99% are not exist.

Table 1 Missingness Analysis

Missing Percentage	Number of Variables	% of All Numerical Variables
100%	0	0%
99% to 100%	0	0%
95% to 99%	0	0%
90% to 95%	0	0%
80% to 90%	0	0%
70% to 80%	0	0%
60% to 70%	0	0%
50% to 60%	0	0%
30% to 50%	0	0%
10% to 30%	1	20%
Below10%	4	80%

The highly positive skewness (with skewness > 10) numerical variables are not exist.

The highly negative skewness (with skewness < -10) numerical variables are not exist.

Table 2 Skewness of Numerical Variables

Skewness Range	Number of Variables	% of All Numerical
----------------	---------------------	--------------------

Table 2 Skewness of Numerical Variables

Skewness Range	Number of Variables	% of All Numerical Variables
10+	0	0%
5 to 10	0	0%
2 to 5	3	60%
1 to 2	0	0%
-1 to 1	2	40%
-2 to -1	0	0%
-5 to -2	0	0%
-10 to -5	0	0%
-10-	0	0%
Total	5	100%

All categorical variables with cardinality over 512 are Name,Ticket.

The calculation of cardinality includes missing category.

The following categorical variables have cell frequency less than 100:

Name,TicketSurvived,Pclass,Sex,Embarked.

CONTENTS

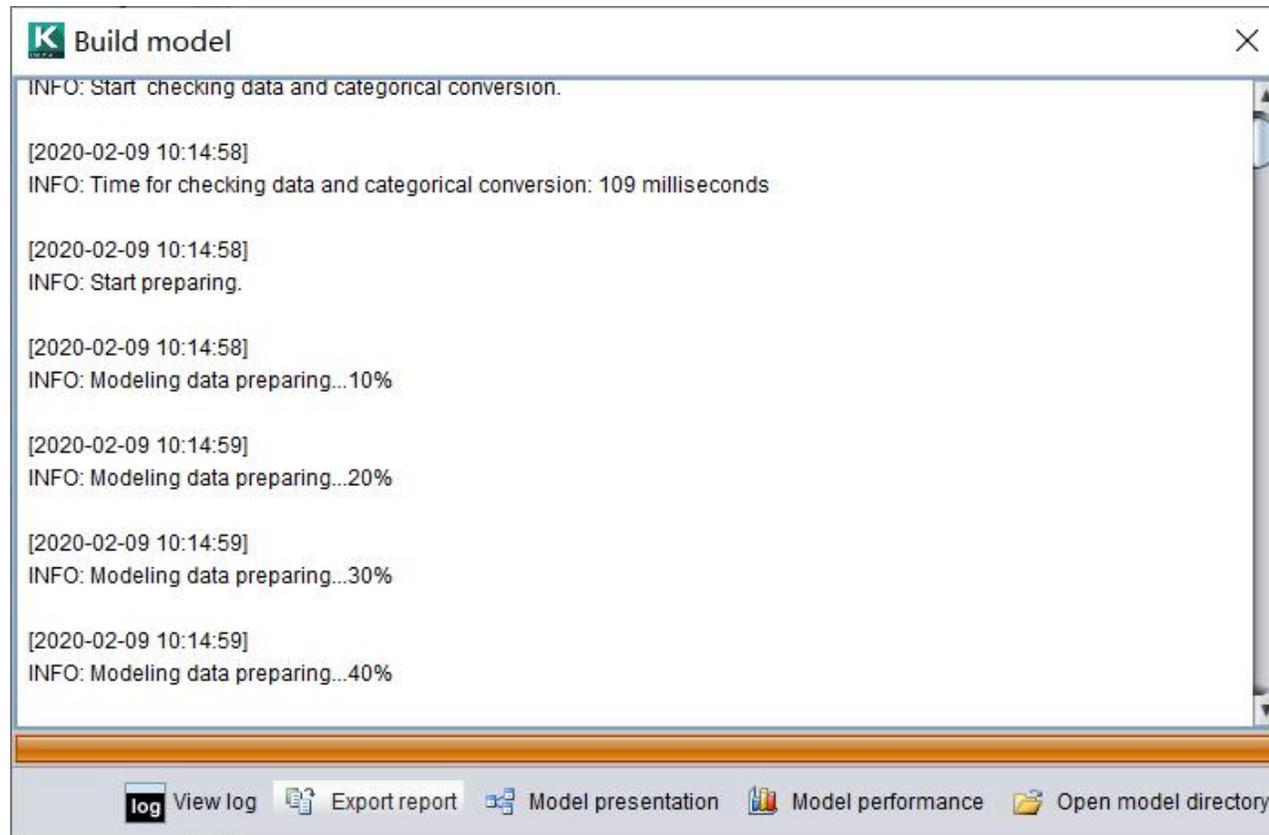
1. Automatic preprocessing
2. Preprocessing report
3. Preprocessing process
4. Manual preprocessing

Preprocessing

➤ 1. Automatic preprocessing



The preprocessing process of intelligent modeling is integrated in the modeling process, with one key automatic preprocessing.



➤ 2. Preprocessing report



After modeling, you can export the model report, which describes the actions of preprocessing. Some contents are as follows:

Target variable: Survived, ID variable: PassengerId.

The number of fields before pretreatment: 12, the number of fields after pretreatment: 11.

The number of fields with missing values before pretreatment: 3 and the number of fields with missing values after pretreatment: 0.

Total rows of data: 891, where deleted rows due to missing target: 0.

Variable selection table			
	Number of selections	Number not selected	Total number
All variables	11	1	12
Unary variables	0	0	0
Binary variables	2	0	2
Category variables	4	1	5
Numerical variables	2	0	2
Counting variables	2	0	2
Datetime variables	0	0	0

Variables Processing Information

Variable name: PassengerId. The type is ID

Variable name: Pclass. The type is Category variables

Number of categories: 3

The variable fills the missing value by using the yimming intelligent filling algorithm.

There are 3 categories are merged because of low frequency.

Generation Category Derivative Variables: BI_Pclass_1, BI_Pclass_2

Variable name: Sex. The type is Binary variables

Number of categories: 2

The variable fills the missing value by using the yimming intelligent filling algorithm.

There are 2 categories are merged because of low frequency.

Generation Category Derivative Variables: BI_Sex_1

Variable name: Age. The type is Numerical variables

Skewness: 0 Average:29.699

Median:24 Variance:13.002

The variable fills the missing value by using the yimming intelligent filling algorithm.

Variable name: SibSp. The type is Counting variables

Skewness: 0 Average:0.523

Median:0 Variance:1.103

The variable fills the missing value by using the yimming intelligent filling algorithm.

Variable name: Parch. The type is Counting variables

➤ 3. Preprocessing process



(1) Check variable value field

Check and record the value range of all variables. If the test data has a category that is not listed in the training data or beyond the range of values, certain processing needs to be carried out.

(2) Time date variable processing

Check all time and date variables and create several commonly used derived variables. Check the correlation of time and date variables, and create multi date linkage derived variables.

(3) Missing value information extraction

If there are missing values in the data, the missing value pattern is extracted and recorded, and the behavior characteristics of missing values are transformed into derivative variables for use.

➤ 3. Preprocessing process



(4) Missing value filling

If there are missing values in the data, use simple or personalized intelligent algorithm to fill in the missing values.

(5) Noise reduction of categorical variables

For the noise that may exist in the categorical variables, such as very few category, abnormal category, suspected error classification and so on, carry out targeted processing.

(6) Convert the categorical variable to a numeric variable

Convert the categorical variable to a numeric variable that can be operated normally. The main method is dummy variable and smoothing, which is judged by algorithm intelligence.

➤ 3. Preprocessing process



(7) Rectify deviation

For some models with normal hypothesis, the high skewness variables are transformed mathematically to make the skewness return to 0, which satisfies the model hypothesis.

(8) Exception handling

Detect and identify possible outliers, and deal with them accordingly.

(9) Variable selection

In order to reduce the time cost and the complexity of the model, we need to remove the useless variables.

➤ 3. Preprocessing process



(10) Standardization / normalization

Data standardization / normalization to eliminate caliber difference. It is beneficial to the optimization of neural networks and other models.

(11) Sample balancing

For binary data, if the proportion of positive and negative samples is seriously unbalanced, it will be balanced according to the specified proportion, and intelligent resampling modeling will be carried out.

➤ 4. Manual preprocessing

Variable selection

Remove some irrelevant variables according to the variable type. For example, ID and long text, single value variable without missing value, etc.

Filter variables according to the importance of variables, only the variables with higher importance are retained. Variable importance can be imported from data dictionary or obtained through modeling.

The screenshot shows the 'Variable filter' dialog box with the 'Importance' filter type selected. The 'Filter type' section has two radio buttons: 'Importance' (selected) and 'Variable type'. Below this, there are three options: 'Select top N by importance' with a value of 3, 'Select variables whose importance is greater than' with a value of 0.357, and 'Only filter selected variables' which is checked. The 'OK' and 'Cancel' buttons are at the bottom right.

The screenshot shows the 'Variable filter' dialog box with the 'Variable type' filter type selected. The 'Filter type' section has two radio buttons: 'Importance' and 'Variable type' (selected). Below this, there is a section titled 'Select by variable type' with several checkboxes: 'Numerical variable' (checked), 'Unary Variable' (unchecked), 'Binary variable' (checked), 'Count variable' (checked), 'Categorical variable' (checked), 'ID' (checked), 'Time and date' (unchecked), and 'Text String' (unchecked). The 'Only filter selected variables' checkbox is also checked. The 'OK' and 'Cancel' buttons are at the bottom right.

4. Manual preprocessing

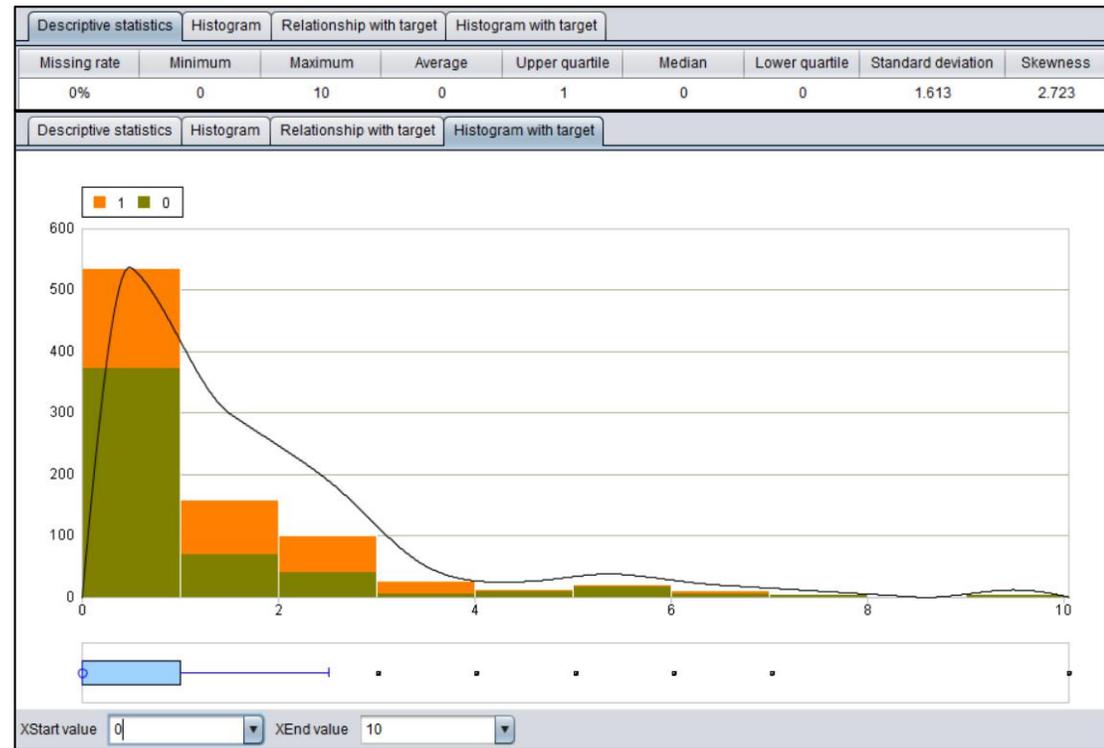


Derived variables

The number of family members is obtained by adding the number of variable "SibSp" and the number of variable "Parch". It can be seen that the survival rate of family members is higher at 1-3.

Variable name	Statistical method	Statistical value
Pclass	Missing rate	0%
Name	Minimum	0
Sex	Maximum	6
Age	Average	0
SibSp	Upper quartile	0
Parch	Median	0
Ticket	Lower quartile	0
Fare	Standard deviation	0.806
Cabin	Skewness	2.744
Embarked		
Family		

Add derived variable family



Variable family statistics

4. Manual preprocessing



Derived variables

The numerical variables can be discretized and converted into categorical variables. Taking age as an example, it is divided into 0, 8, 18, 35 and 60 age groups, generating derivative variables and making statistics.

Computed variable name: AgeArea

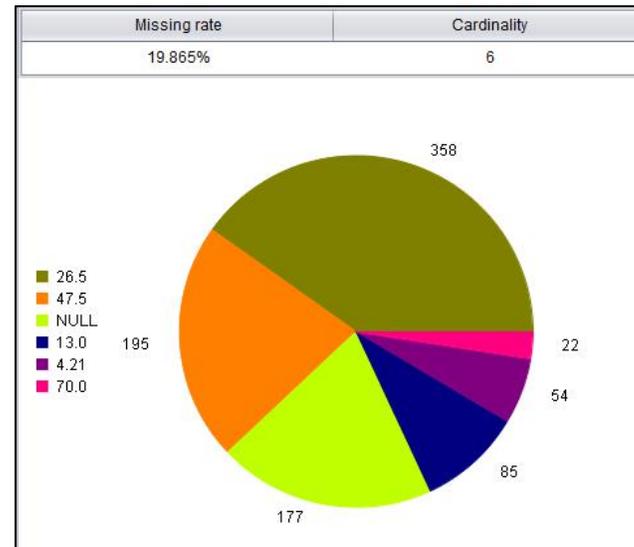
Variable name: Age

Bin boundaries:

NO.	Bin boundaries
1	Minimum: 0.42
2	8
3	18
4	35
5	60
6	Maximum: 80.0

Missing rate: 19.865%

Add derived variable AgeArea



Categorical Level	Frequency	Positive Frequency	Positive Ratio	Odds wrt All
4.21	54	36	66.667%	1.737
47.5	195	78	40%	1.042
13.0	85	34	40%	1.042
26.5	358	137	38.268%	0.997
NULL	177	52	29.379%	0.765
70.0	22	5	22.727%	0.592
All	891	342	38.384%	1

Variable AgeArea statistics

It can be seen that the survival rate of the 0-8-year-old is the highest, the difference between the young and the middle-aged is not big, and the survival rate of the old is the lowest.

➤ 4. Manual preprocessing

Preprocessing options

In the model options, you can define whether to preprocess data and whether to fill it intelligently.

If the data has been preprocessed, you can cancel the data preprocessing.

Intelligent filling can better fill the missing value, but it will consume more hardware resources and time. When the amount of data is large, intelligent filling is not recommended. If unchecked, it will be filled in simply.

The screenshot shows the 'Model options' dialog box with the following settings:

- Data preparation
- Intelligent impute
- Resampling: Number of samples: 5
- Best number of sample combinations: 3
- Balanced sampling ratio: 1:1
- Sample multiplier: 150
- Ensemble method: Optimal model strategy
- Best number of ensembles: 0
- Ensemble function: np.mean
- Model evaluation criterion: (empty)
- Percentage of test data: Automatic %
- Adjust scoring results
- Set random seeds: 0

Buttons: OK, Cancel

CONTENTS

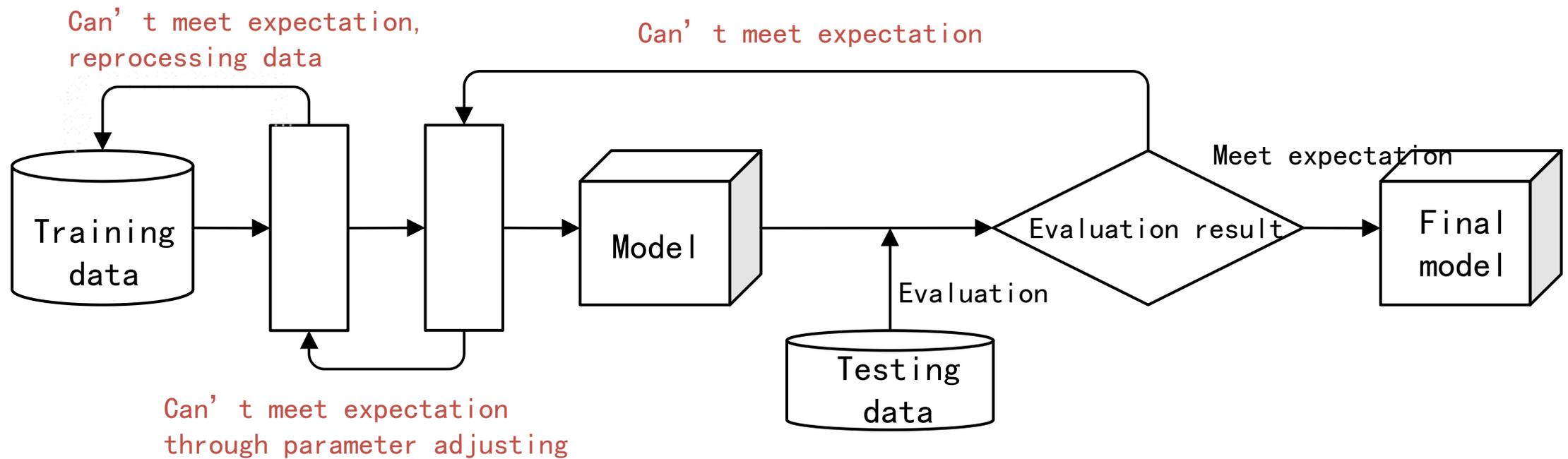
1. Modeling process
2. Intelligent modeling
3. Professional modeling

04
Modeling

1. Modeling process



When using traditional tools, it usually requires professionals with statistical basis to continuously select algorithms, adjust model parameters, and finally get the expected model. The modeling process is as follows:



2. Intelligent modeling



Intelligent modeling tools do not need statistical knowledge, one key intelligent modeling, optimization of model combination and model parameters are implemented internally.

```
2020-02-09 10:15:20,085 - yiming_model.cp37-win_amd64.pyd[line:90] - INFO: Model built successfully
2020-02-09 10:15:20,085 - yiming_model.cp37-win_amd64.pyd[line:90] - DEBUG: feature importance of YiModel: {'Rank_F
are': 1.0, 'Pow0_69_Age': 0.6549645954182951, 'MI_Age': 0.43832267557855187, 'Rank_SibSp': 0.39012433562963306, '
Rank_Parch': 0.0}
2020-02-09 10:15:20,085 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: performance of each base model in Yi
Model: {'XGBClassification_1': 0.8122683142100618, 'RFClassification_1': 0.7518387761106208, 'FNNClassification_1': 0
.5, 'RidgeClassification_1': 0.757811120917917, 'TreeClassification_1': 0.7086201824065902, 'LogicClassification_1': 0.7
496322447778758, 'CNNClassification_1': 0.5, 'GBDTClassification_1': 0.7994998528979111}
2020-02-09 10:15:20,085 - interface_library.cp37-win_amd64.pyd[line:90] - INFO: Calculate predict value on test data
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: predict value on test data:
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - INFO: Calculate ensemble performance
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: ensemble performance: 0.820535
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - INFO: Writing out results
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: writing out predict values
2020-02-09 10:15:20,132 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: writing out model
2020-02-09 10:15:20,155 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: writing out feature importance
2020-02-09 10:15:20,155 - interface_library.cp37-win_amd64.pyd[line:90] - DEBUG: writing out modeling information
2020-02-09 10:15:20,155 - interface_library.cp37-win_amd64.pyd[line:90] - INFO: Build model finished
```

log View log Export report Model presentation Model performance Open model directory

➤ 3. Professional modeling



Intelligent modeling opens up model parameters for professional users who are proficient in the models. Here are the general options for the model:

K Model options [Close]

Normal | Binary model | Regression model | Multiclassification model

Data preparation Intelligent impute

Resampling Number of samples: 5 Best number of sample combinations: 3

Balanced sampling ratio: 1:1 Sample multiplier: 150

Ensemble method: Optimal model strategy Best number of ensembles: 0

Ensemble function: np.mean Model evaluation criterion: [Dropdown]

Percentage of test data: Automatic %

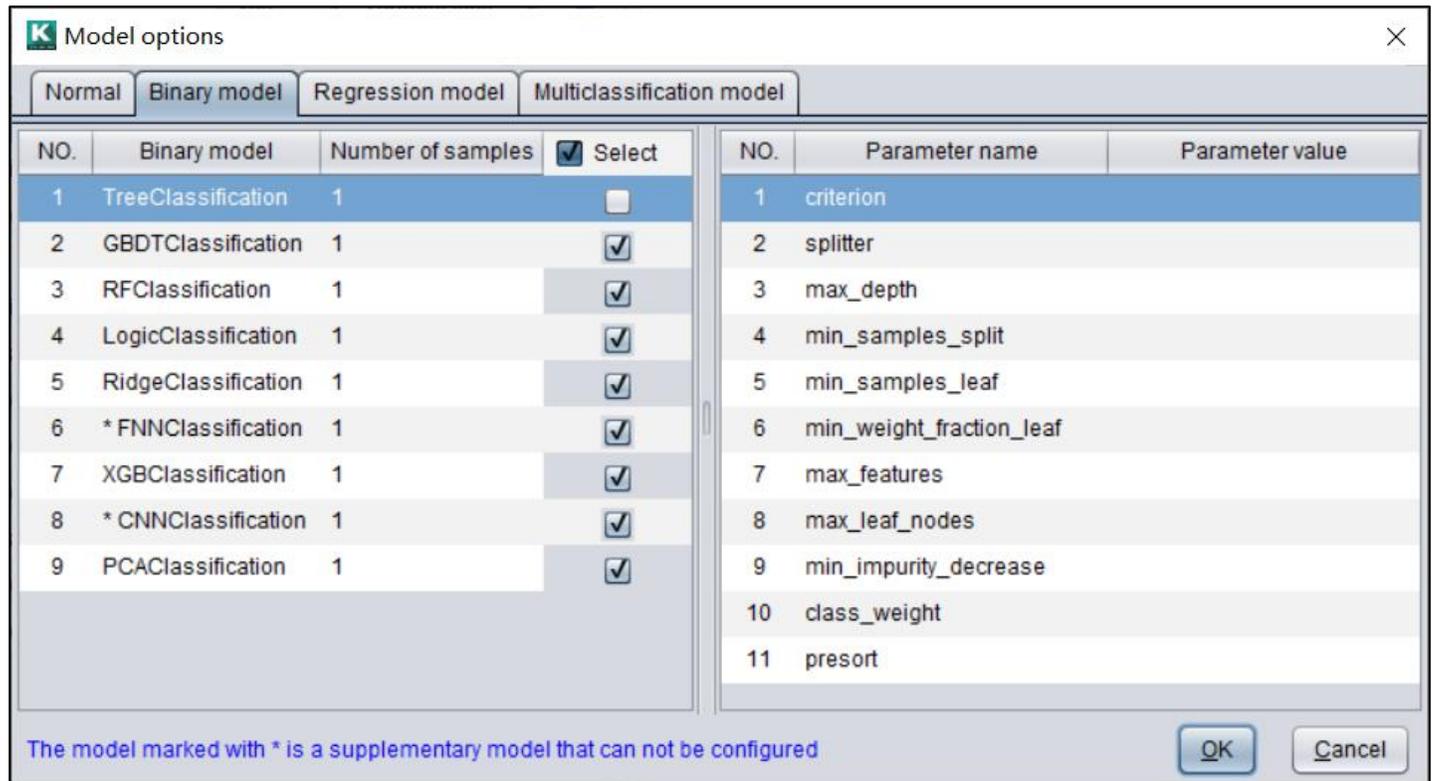
Adjust scoring results Set random seeds: 0

[OK] [Cancel]

➤ 3. Professional modeling



Intelligent modeling supports several binary classification algorithm models in the graph, and can also set whether each model is used and the sampling times. On the right, you can set parameter values for each model. For ordinary users, these settings can be ignored.



➤ 3. Professional modeling



Similarly, we can set whether to use regression model and multi classification model, and their respective parameters.

NO.	Regression model	Number of samples	Select	NO.	Parameter name	Parameter value
1	TreeRegression	1	<input type="checkbox"/>	1	criterion	
2	GBDTRegression	1	<input checked="" type="checkbox"/>	2	splitter	
3	RFRRegression	1	<input checked="" type="checkbox"/>	3	max_depth	
4	LRegression	1	<input type="checkbox"/>	4	min_samples_split	
5	LassoRegression	1	<input checked="" type="checkbox"/>	5	min_samples_leaf	
6	ENRegression	1	<input checked="" type="checkbox"/>	6	min_weight_fraction_leaf	
7	RidgeRegression	1	<input checked="" type="checkbox"/>	7	max_features	
8	* FNNRegression	1	<input checked="" type="checkbox"/>	8	max_leaf_nodes	
9	XGBRegression	1	<input checked="" type="checkbox"/>	9	min_impurity_decrease	
10	* CNNRegression	1	<input checked="" type="checkbox"/>	10	presort	
11	PCARegression	1	<input checked="" type="checkbox"/>			

The model marked with * is a supplementary model that can not be configured

NO.	Multiclassification model	Number of samples	Select	NO.	Parameter name	Parameter value
1	XGBMultiClassification	1	<input checked="" type="checkbox"/>	1	max_depth	
2	* CNNMultiClassification	1	<input checked="" type="checkbox"/>	2	learning_rate	
				3	n_estimators	
				4	booster	
				5	gamma	
				6	min_child_weight	
				7	max_delta_step	
				8	subsample	
				9	colsample_bytree	
				10	colsample_bylevel	
				11	reg_alpha	
				12	reg_lambda	

The model marked with * is a supplementary model that can not be configured

Detailed documentation of each model parameter : <http://doc.raqsoft.com.cn/AIModel/userrefer/jm20.html>

CONTENTS

1. Model performance
2. Model presentation
3. Variable importance

05
Model performance



➤ 1. Model performance

Classification model: evaluation index

Intelligent modeling provides three commonly used evaluation indexes for classification model:

GINI	AUC	KS
0.785054	0.892527	0.670079

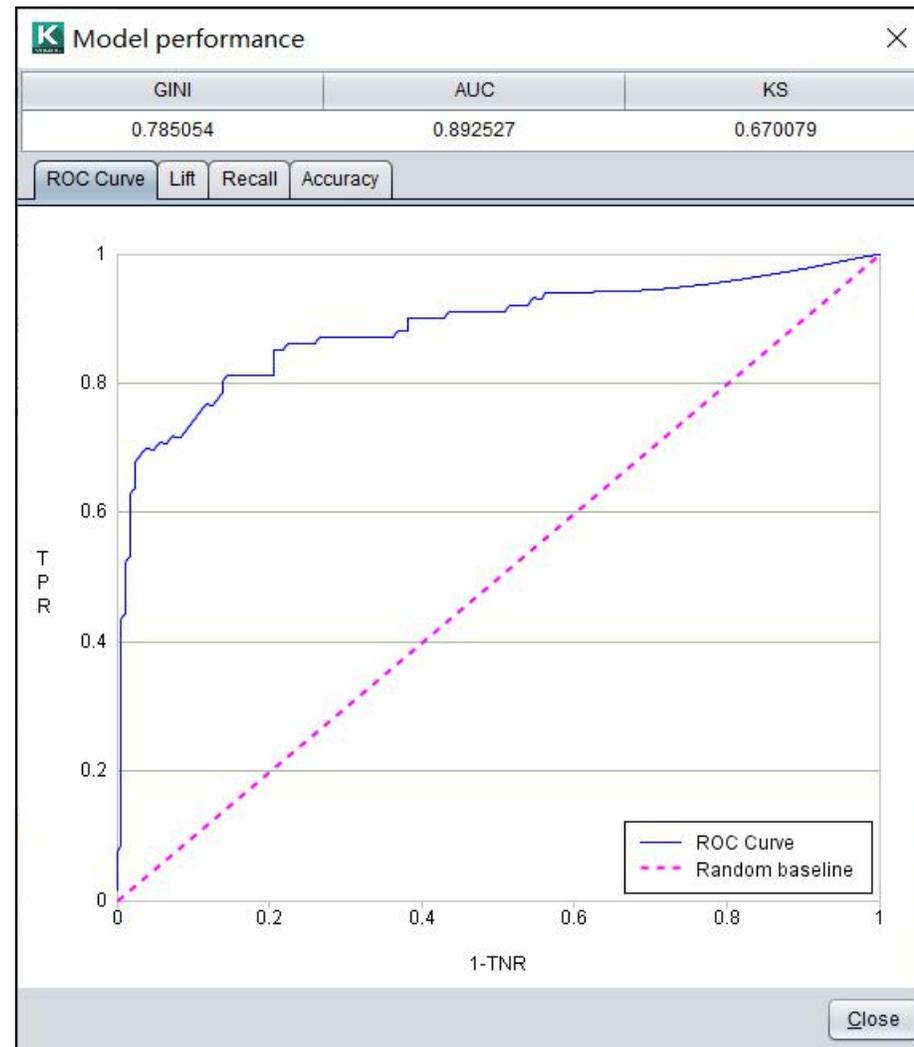
Evaluation Index	Description
GINI	Gini index is equal to $2 * auc - 1$ in numerical value, which is used to characterize the model's ability to distinguish positive and negative samples.
AUC	AUC is equal to the area under ROC curve. The higher AUC is, the better the model is.
KS	KS value is used to measure the ability of the model to distinguish positive and negative samples. The larger the KS value is, the stronger the ability of the model to distinguish positive and negative samples is.

➤ 1. Model performance



Classification model: ROC curve

ROC curve is the relationship between true positive class rate and "1-true negative class rate". ROC curve can be regarded as a visual display to evaluate all possible decision-making performance of a given model.



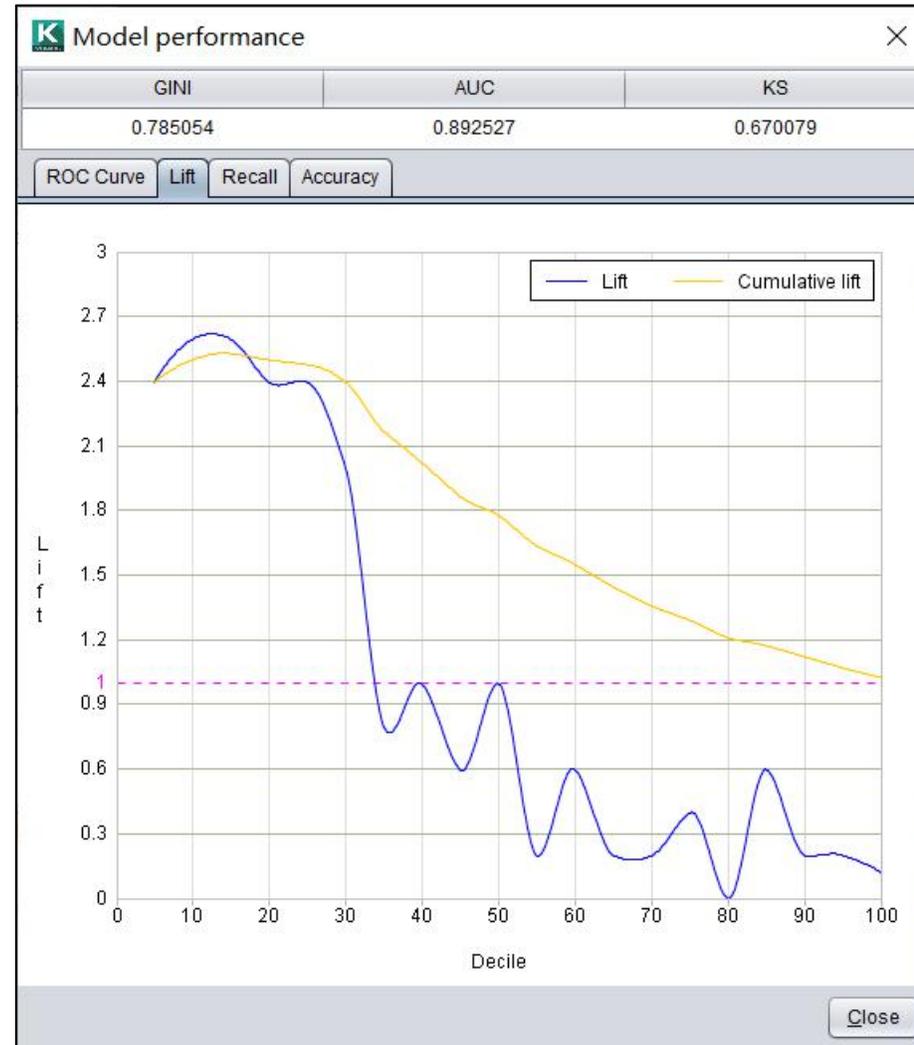
➤ 1. Model performance



Classification model: Lift

Lift refers to the multiple that can be improved by using association rules. It is the ratio of the degree of confidence to expected confidence.

Lift is particularly suitable for targeted marketing and other scenarios.

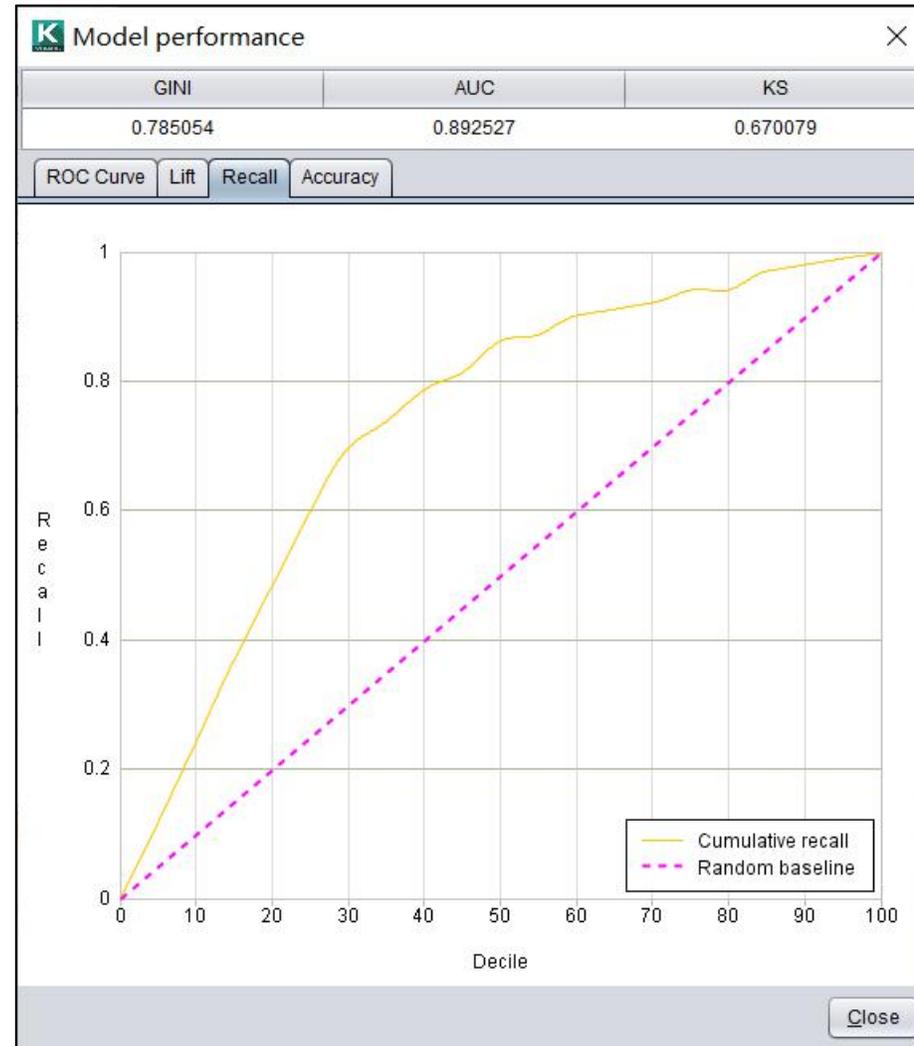


➤ 1. Model performance



Classification model: Recall

Recall graph shows that the model can find positive samples, which is mainly used in the scene of data imbalance. The cumulative recall rate is the ratio of cumulative positive samples and total positive samples in each group.



➤ 1. Model performance

Classification model: Accuracy table

Threshold: value used to distinguish positive and negative samples.

Accuracy: the ratio of correct samples to all samples.

Precision: the correct rate of prediction in the result of positive sample.

Recall: the ratio of correctly predicted positive samples and all positive samples.

Model performance window showing GINI, AUC, and KS values, and an Accuracy table with columns for Threshold, Accuracy, Precision, and Recall.

GINI	AUC	KS
0.785054	0.892527	0.670079

ROC Curve | Lift | Recall | Accuracy

Lower limit: 0.05 | Upper limit: 0.95 | Number of subsections: 19 | Set

Threshold	Accuracy	Precision	Recall
0.05	0.448	0.41	0.99
0.1	0.552	0.46	0.942
0.15	0.653	0.528	0.913
0.2	0.728	0.599	0.883
0.25	0.787	0.674	0.864
0.3	0.813	0.715	0.854
0.35	0.806	0.714	0.825
0.4	0.836	0.776	0.806
0.45	0.828	0.782	0.767
0.5	0.843	0.843	0.728
0.55	0.854	0.9	0.699
0.6	0.854	0.944	0.66
0.65	0.847	0.943	0.641
0.7	0.84	0.955	0.612
0.75	0.81	0.964	0.524
0.8	0.799	0.962	0.495
0.85	0.776	0.957	0.437

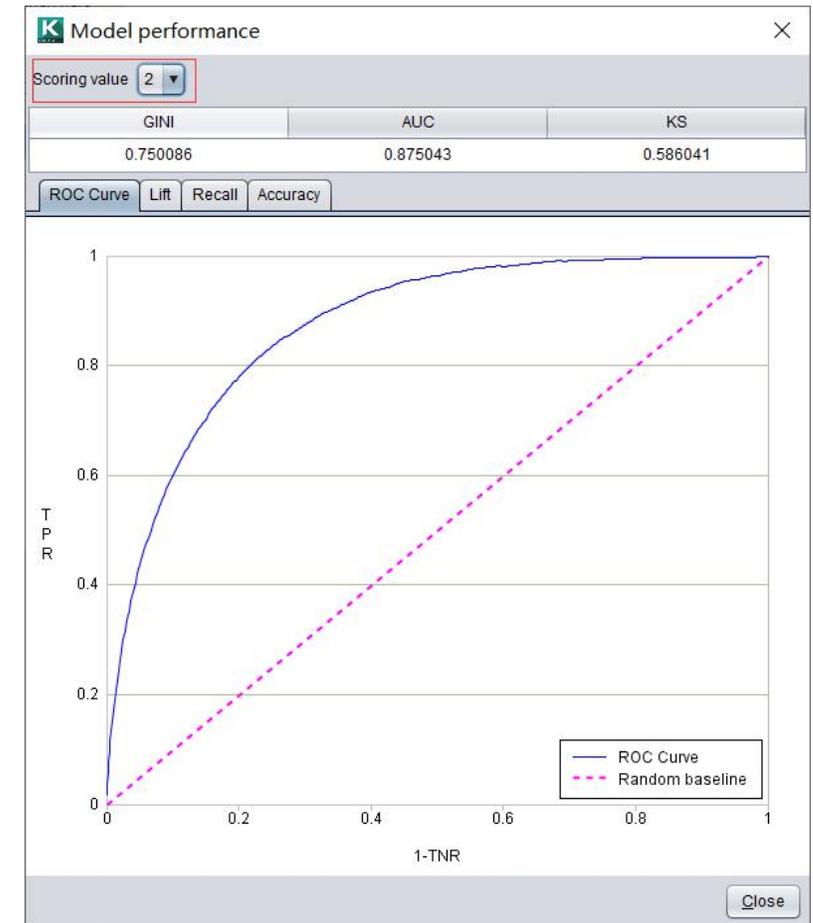
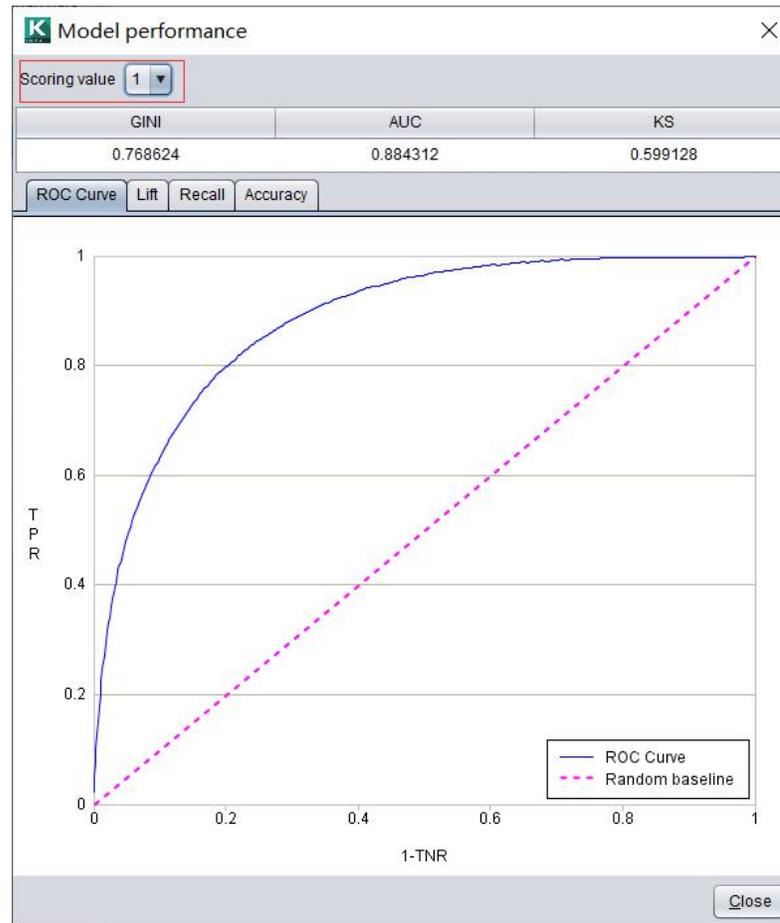
Close

➤ 1. Model performance



Multiclassification model

When the target variable is a categorical variable, the model performance of each classification can be viewed by switching prediction values.

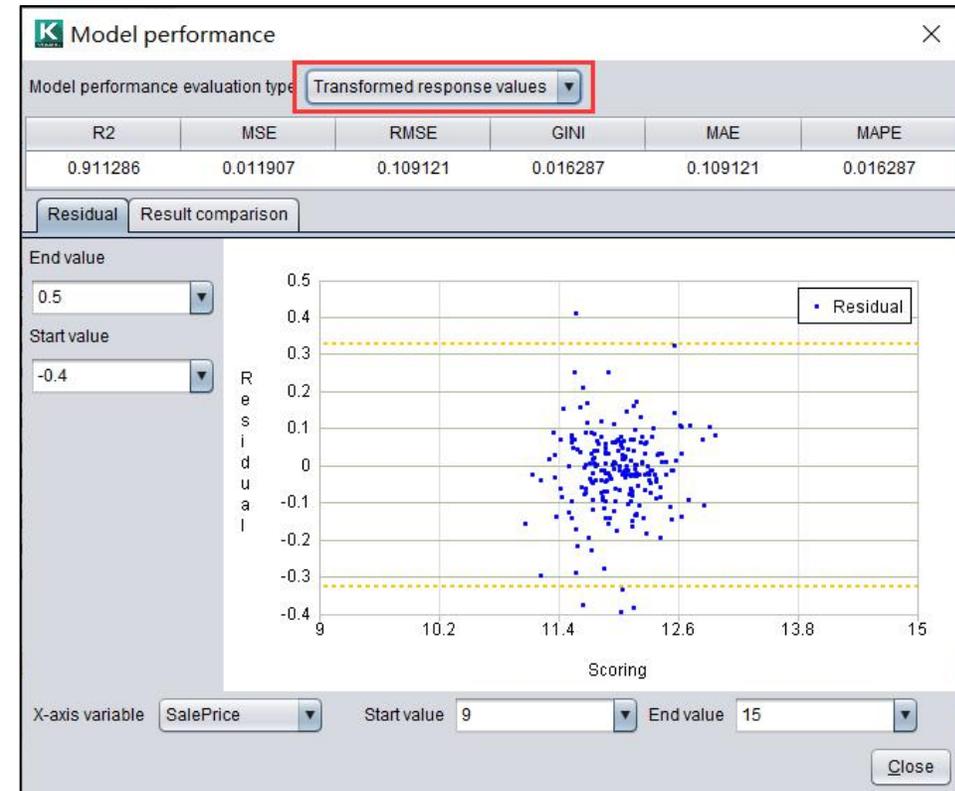
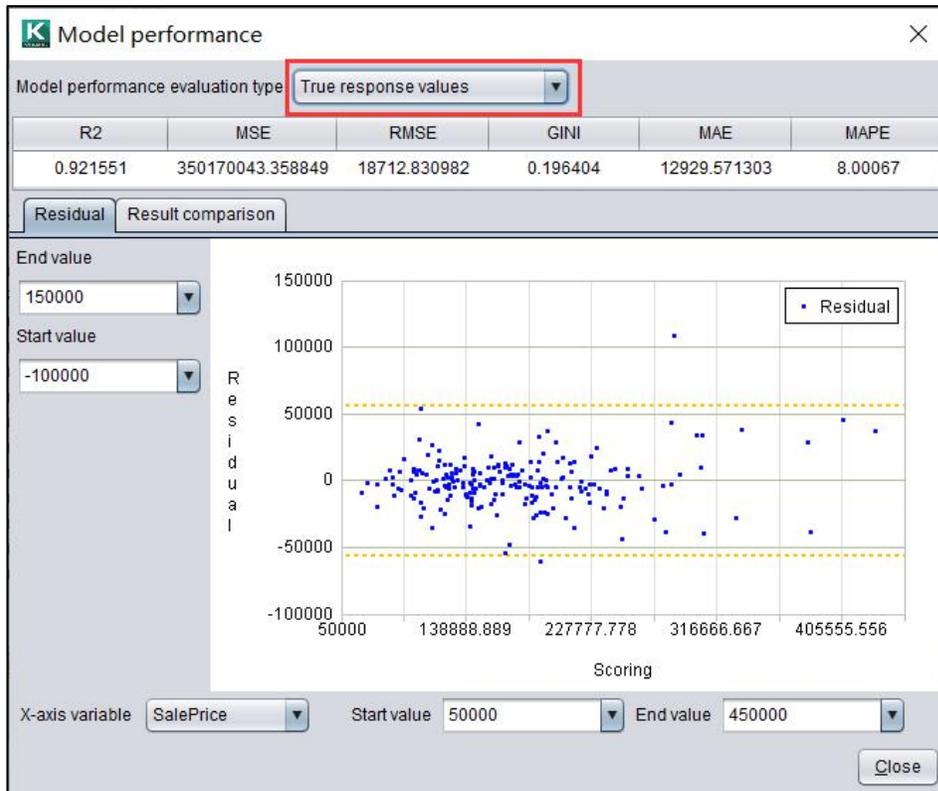


➤ 1. Model performance



Regression model: True response values and transformed response values

The performance of regression model can be divided into true value performance and transformed value performance (data value after preprocessed). The true value looks more intuitive, and the transformed value is more accurate for the evaluation of model performance.



➤ 1. Model performance



Regression model: evaluation index

Intelligent modeling provides six commonly used evaluation indexes of regression model:

Model performance evaluation type: True response values					
R2	MSE	RMSE	GINI	MAE	MAPE
0.921551	350170043.358...	18712.830982	0.196404	12929.571303	8.00067

Model performance evaluation type: Transformed response values					
R2	MSE	RMSE	GINI	MAE	MAPE
0.911286	0.011907	0.109121	0.016287	0.109121	0.016287

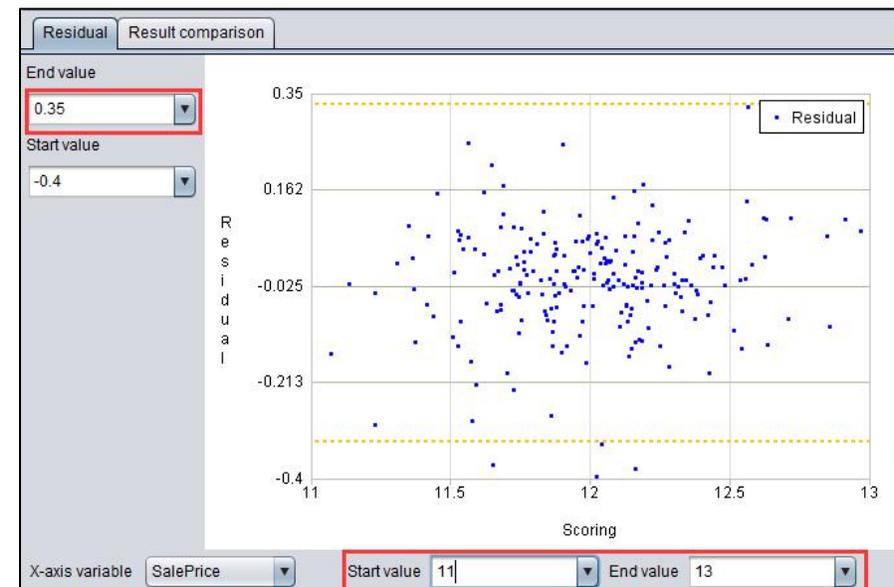
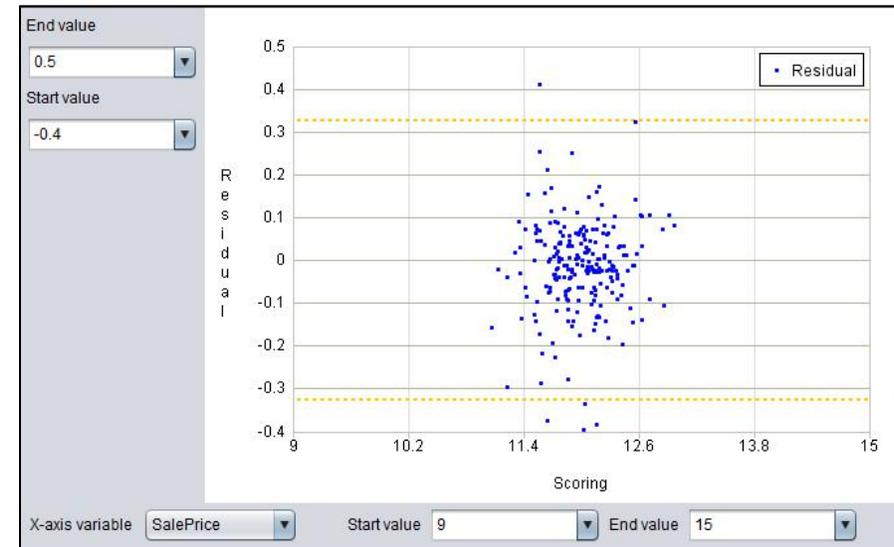
Evaluation Index	Description
R ²	R ² is the ratio of the sum of the square of the error between the predicted value and the observed value to the sum of the square of the difference between the observed value and the observed mean value.
MSE	The average sum of the squares of the deviations of the predicted value from the true value.
RMSE	The square root of MSE. The order of magnitude is the same as the true value.
GINI	The average of the absolute value of the deviation between the predicted value and the true value.
MAE	The average of the absolute value of the deviation between the predicted value and the true value.
MAPE	The average of the absolute value of the deviation between the predicted value and the true value.

➤ 1. Model performance

Regression model: residual chart

The residual is the difference between the observed value and the predicted value. The residual chart is a scatter chart with the residual as the vertical axis and any numerical variable as the horizontal axis. The yellow line is three times RMSE.

You can adjust the horizontal axis variable and the value range of the horizontal and vertical axis for further viewing.



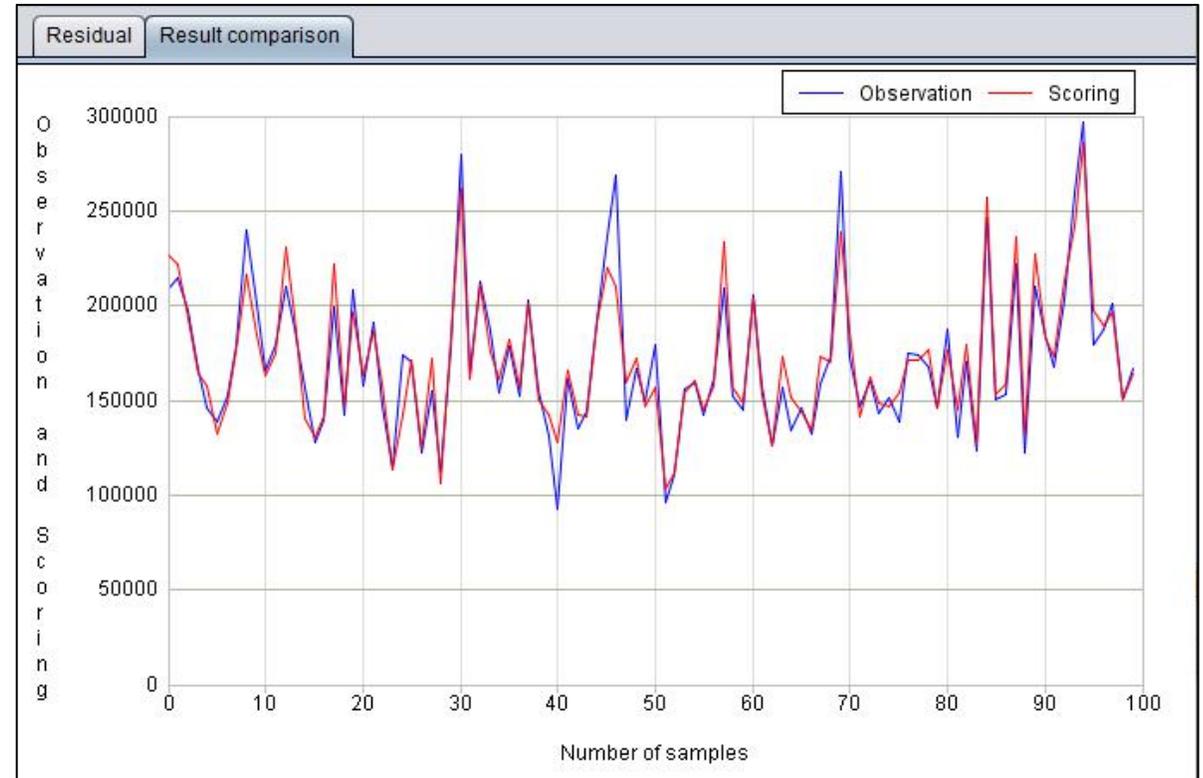
➤ 1. Model performance



Regression model: result comparison chart

The horizontal axis of the result comparison chart are the samples of random distribution, and the vertical axis is the corresponding observation value and prediction value.

Blue is the observed value and red is the predicted value.



2. Model presentation



The model presentation lists the final selected model combinations and the parameter values of each model. The selected model parameters can be copied to the model options through the button to further optimize the model parameters.

Ensemble performance: 0.892527

Model name	auc	Select
XGBClassification_1	0.879464	<input checked="" type="checkbox"/>
* FNNClassification_1	0.873433	<input checked="" type="checkbox"/>
RidgeClassification_1	0.872050	<input checked="" type="checkbox"/>
GBDTClassification_1	0.882200	<input checked="" type="checkbox"/>

Unused models	auc	Select
RFClassification_1	0.846454	<input type="checkbox"/>
LogicClassification_1	0.865166	<input type="checkbox"/>
* CNNClassification_1	0.786437	<input type="checkbox"/>
PCAClassification_1	0.851927	<input type="checkbox"/>

Parameter name	Parameter value
max_depth	6
learning_rate	0.1
n_estimators	150
objective	binary:logistic
booster	gbtree
gamma	0
min_child_weight	1
max_delta_step	0
subsample	1
colsample_bytree	1
colsample_bylevel	1
reg_alpha	0
reg_lambda	1
scale_pos_weight	1

The model marked with * is a supplementary model that can not be configured

Copy selected model to model options

Close

The final classification model and parameters of Titanic model

Ensemble performance: 296270797.147141

Model name	mse	Select
GBDTRegression_1	382347158.307895	<input checked="" type="checkbox"/>
LassoRegression_1	445564549.912437	<input checked="" type="checkbox"/>
XGBRegression_1	369838540.587562	<input checked="" type="checkbox"/>

Unused models	mse	Select
ENRegression_1	445674473.7744...	<input type="checkbox"/>
PCARegression_1	596129767.9753...	<input type="checkbox"/>
* CNNRegression_1	17585172054.03...	<input type="checkbox"/>
RidgeRegression_1	551359413.2697...	<input type="checkbox"/>
RFRRegression_1	812798684.5407...	<input type="checkbox"/>
* FNNRegression_1	659157565.9669...	<input type="checkbox"/>

Parameter name	Parameter value
loss	ls
learning_rate	0.1
n_estimators	100
subsample	1.0
criterion	friedman_mse
min_samples_split	50
min_samples_leaf	50
min_weight_fraction_leaf	0
max_depth	6
min_impurity_decrease	1e-08
max_features	null
alpha	0.9
max_leaf_nodes	null
warm_start	false
presort	

The model marked with * is a supplementary model that can not be configured

Copy selected model to model options

Close

The final regression model and parameters of house price model

➤ 3. Variable importance



After modeling, the importance information of each variable can be obtained. From the returned importance of Titanic model, we can see that age (children first) and ticket price (higher class) are the most important factors for survival.

NO.	Variable name	Type	Date format	<input checked="" type="checkbox"/> Select	Importance
1	Sex	Binary variable		<input checked="" type="checkbox"/>	1
2	AgeArea	Categorical variable		<input checked="" type="checkbox"/>	0.726
3	Pclass	Categorical variable		<input checked="" type="checkbox"/>	0.524
4	SibSp	Categorical variable		<input checked="" type="checkbox"/>	0.443
5	Age	Numerical variable		<input checked="" type="checkbox"/>	0.392
6	Fare	Numerical variable		<input checked="" type="checkbox"/>	0.275
7	Parch	Categorical variable		<input checked="" type="checkbox"/>	0.244
8	Family	Numerical variable		<input checked="" type="checkbox"/>	0.197
9	Cabin	Categorical variable		<input checked="" type="checkbox"/>	0.169
10	Embarked	Categorical variable		<input checked="" type="checkbox"/>	0.146
11	PassengerId	ID		<input checked="" type="checkbox"/>	0
12	Survived	Binary variable		<input checked="" type="checkbox"/>	-
13	Name	ID		<input type="checkbox"/>	0
14	Ticket	Categorical variable		<input checked="" type="checkbox"/>	0

The role of variable importance	
1	Refer to the importance of variables and reprocess the data accordingly.
2	The important variables are used interactively to generate the derived variables, such as distance / time = speed, speed * time = distance and so on.
3	Refer to the importance of variables and make targeted suggestions to customers.

CONTENTS

1. Batch prediction
2. Single prediction



Prediction

➤ 1. Batch prediction



After you create the model, you can use test data for prediction.

For the binary classification model, the first column is the probability that the target variable is a positive sample.

Taking Titanic as an example, the probability of survival of No. 624 passenger is predicted to be 32.984%.

Survived_1_percentage	PassengerId	Survived	Pclass	Name	Sex
21.584%	624	0	3	Hansen, Mr. Henry Damsgaard	male
13.652%	625	0	3	"Bowen, Mr. David John ""Dai""	male
21.625%	626	0	1	Sutton, Mr. Frederick	male
9.799%	627	0	2	Kirkland, Rev. Charles Leonard	male
95.103%	628	1	1	Longley, Miss. Gretchen Fiske	female
12.653%	629	0	3	Bostandyeff, Mr. Guentcho	male
6.248%	630	0	3	O'Connell, Mr. Patrick D	male
47.066%	631	1	1	Barkworth, Mr. Algernon Henry Wilson	male
3.796%	632	0	3	Lundahl, Mr. Johan Svensson	male
63.63%	633	1	1	Stahelin-Maeglin, Dr. Max	male
7.895%	634	0	1	Parr, Mr. William Henry Marsh	male
17.128%	635	0	3	Skoog, Miss. Mabel	female
87.152%	636	1	2	Davis, Miss. Mary	female
30.58%	637	0	3	Leinonen, Mr. Antti Gustaf	male
26.677%	638	0	2	Collyer, Mr. Harvey	male
33.02%	639	0	3	Panula, Mrs. Juha (Maria Emilia Ojala)	female
9.154%	640	0	3	Thorneycroft, Mr. Percival	male
23.667%	641	0	3	Jensen, Mr. Hans Peder	male
97.429%	642	1	1	Sagesser, Mlle. Emma	female
50.589%	643	0	3	Skoog, Miss. Margit Elizabeth	female
24.772%	644	1	3	Foo, Mr. Choong	male
87.833%	645	1	3	Baclini, Miss. Eugenie	female

➤ 1. Batch prediction



For regression model, the first column is the predicted value of the target variable.

Taking the house price as an example, the price of house 1461 is predicted to be 120644.118.

Batch scoring		Scoring						
Scoring data		C:\Users\wunan\OneDrive\data\house_prices_test.csv						
SalePrice_predictvalue	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape
129298.66	1461	20	RH	80	11622	Pave		Reg
160807.103	1462	20	RL	81	14267	Pave		IR1
191135.414	1463	60	RL	74	13830	Pave		IR1
198392.522	1464	60	RL	78	9978	Pave		IR1
189149.272	1465	120	RL	43	5005	Pave		IR1
168192.263	1466	60	RL	75	10000	Pave		IR1
185509.826	1467	20	RL		7980	Pave		IR1
158992.343	1468	60	RL	63	8402	Pave		IR1
200264.592	1469	20	RL	85	10176	Pave		Reg
120879.556	1470	20	RL	70	8400	Pave		Reg
200707.594	1471	120	RH	26	5858	Pave		IR1
98007.484	1472	160	RM	21	1680	Pave		Reg
97726.594	1473	160	RM	21	1680	Pave		Reg
141444.443	1474	160	RL	24	2280	Pave		Reg
103150.052	1475	120	RL	24	2280	Pave		Reg
354069.211	1476	60	RL	102	12858	Pave		IR1
255800.709	1477	20	RL	94	12883	Pave		IR1
282393.717	1478	20	RL	90	11520	Pave		Reg
313624.796	1479	20	RL	79	14122	Pave		IR1
490093.299	1480	20	RL	110	14300	Pave		Reg
331277.321	1481	60	RL	105	13650	Pave		Reg

➤ 1. Batch prediction



When the target variable is a categorical variable, the probability (sum of 1) of each target classification value is displayed after prediction. For example, for the first record, the probability of target value of 2 is the highest, which is 97.402%.

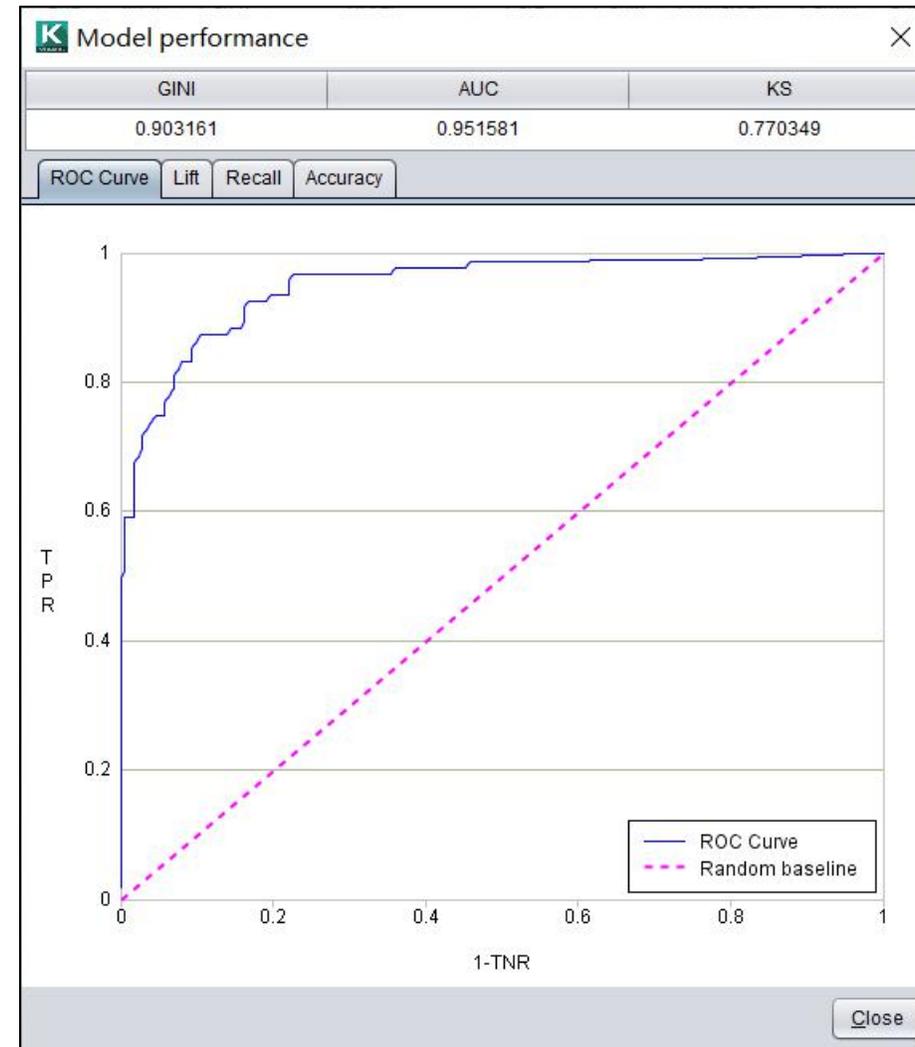
Cover_Type_1_percentage	Cover_Type_2_percentage	Cover_Type_3_percentage	Cover_Type_4_percentage	Cover_Type_5_percentage	Cover_Type_6_percentage	Cover_Type_7_percentage
0.448%	97.402%	0.169%	0.021%	1.745%	0.177%	0.038%
0.297%	98.152%	0.115%	0.015%	1.223%	0.172%	0.027%
1.875%	97.405%	0.594%	0.01%	0.088%	0.011%	0.017%
3.302%	94.912%	1.172%	0.014%	0.146%	0.429%	0.025%
0.319%	97.864%	0.091%	0.014%	1.546%	0.137%	0.027%
0.768%	96.389%	0.337%	0.034%	2.059%	0.359%	0.054%
0.699%	95.365%	0.171%	0.021%	3.529%	0.176%	0.039%
0.37%	96.957%	0.095%	0.015%	2.385%	0.148%	0.029%
0.511%	97.973%	0.107%	0.014%	1.211%	0.163%	0.021%
0.673%	98.115%	0.073%	0.013%	0.999%	0.103%	0.024%
0.421%	98.58%	0.137%	0.011%	0.708%	0.124%	0.019%
3.994%	95.644%	0.044%	0.022%	0.222%	0.026%	0.047%
2.927%	96.683%	0.178%	0.01%	0.155%	0.028%	0.019%
0.229%	98.33%	0.182%	0.011%	1.119%	0.11%	0.018%
0.318%	98.225%	0.133%	0.024%	0.802%	0.448%	0.05%
0.704%	94.935%	0.224%	0.041%	2.922%	1.099%	0.074%
1.336%	96.347%	0.178%	0.033%	1.74%	0.315%	0.052%
0.383%	96.798%	0.146%	0.027%	2.265%	0.329%	0.053%
0.234%	94.256%	0.108%	0.016%	5.058%	0.297%	0.03%
0.252%	96.695%	0.12%	0.019%	2.536%	0.335%	0.043%
0.681%	97.92%	0.139%	0.029%	0.718%	0.452%	0.06%
6.872%	92.693%	0.026%	0.018%	0.334%	0.021%	0.035%

➤ 1. Batch prediction



Generally, the prediction data does not contain the target variable.

When target variable is included in the prediction data, the performance of the model can be calculated according to the prediction result to evaluate the model.

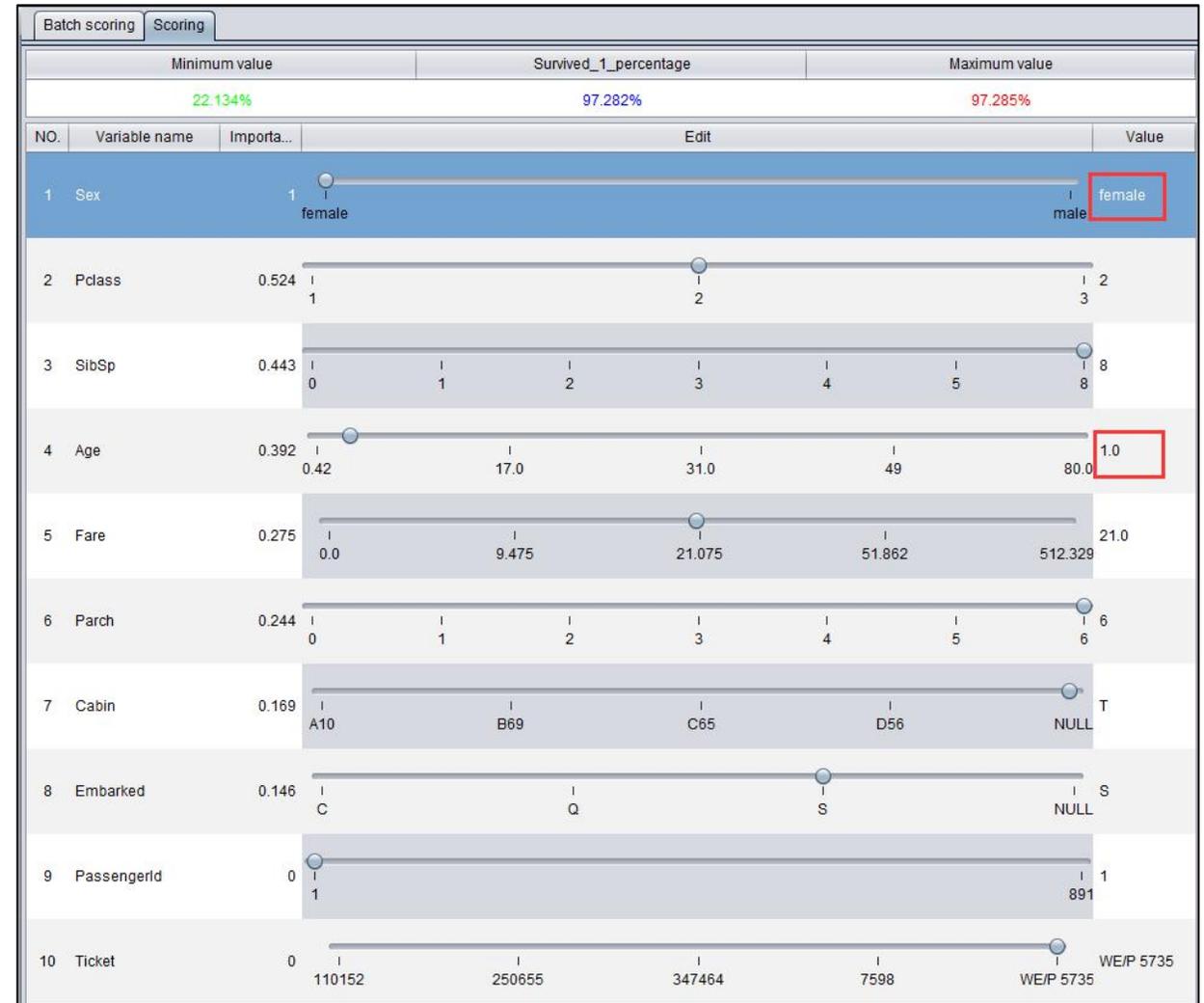


➤ 2. Single prediction



A single prediction can be dragged to modify the variable value and view the prediction result in real time.

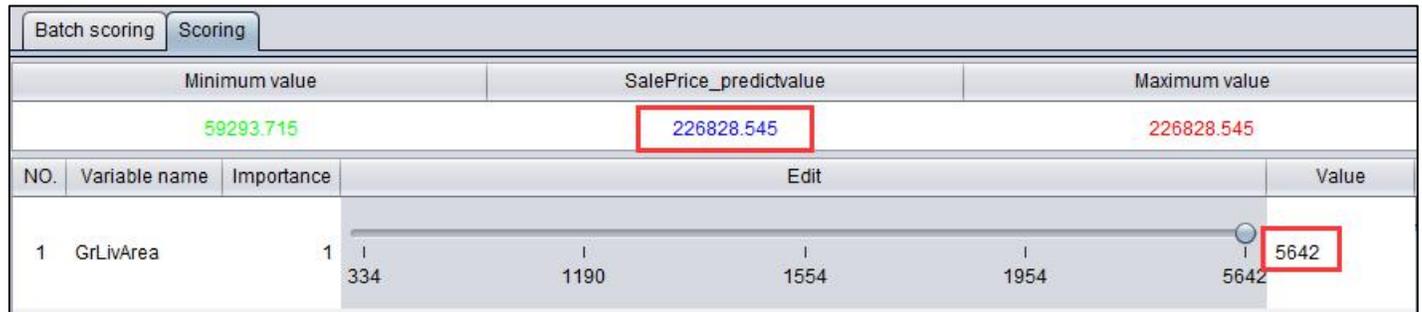
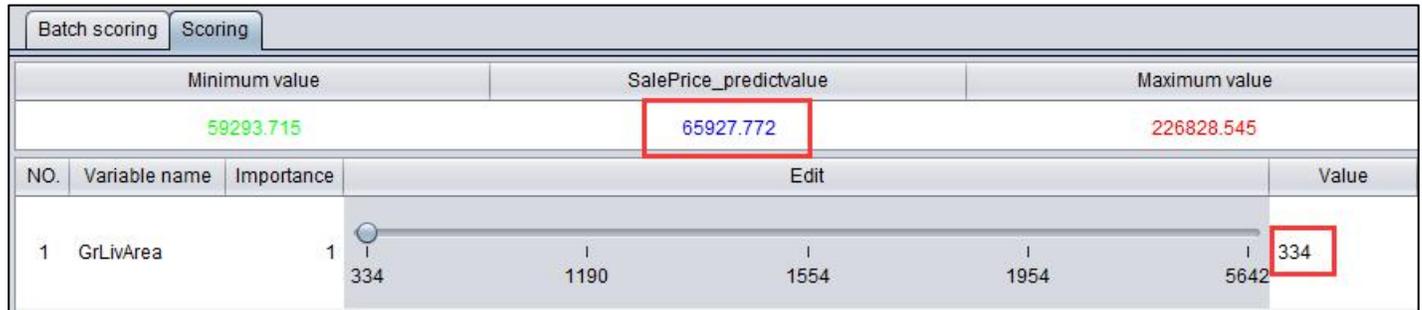
The variables are arranged in descending order of importance, and the top variables usually have more influence on the prediction result. It can be seen that the survival rate of the younger females is very high.



2. Single prediction



For the house price prediction model, we can see that when the basement area is dragged from 334 to 5642 (other variables have not changed), the house price has greatly increased.



CONTENTS

1. esProc External library
2. Integration architecture

07
Integration solution

➤ 1. esProc External library



esProc external library provides interface functions for intelligent modeling, which can be called by SPL. The SPL for modeling:

	A	B
1	=file("titanic_train.csv").cursor@cqt()	/Create training data cursor
2	=ym_env()	/Initialize environment
3	=ym_model(A2,A1)	/Loading data
4	=ym_target(A3, "Survived")	/Set target variable
5	=ym_build_model(A3)	/Execute modeling
6	=ym_save_pcf(A5,"titanic.pcf")	/Save model file
7	=ym_json(A5)	/Export model information as JSON string
8	=ym_importance(A5)	/Get variable importance
9	=ym_present(A5)	/Get model presentation
10	=ym_performance(A5)	/Get model performance
11	>ym_close(A2)	/Close

A7

Value
{"Importance":{"PassengerId":0,"Pclass":0,"Sex":0,"Age":0.433191...

A8

Name	Importance
PassengerId	0.0
Pclass	0.0
...	...

A9

name	value	properties
XGBClass...	0.815	[[max_delt...
XGBClass...	0.777	[[max_delt...
...

A10

Name	Value
GINI	0.617
AUC	0.808
...	...

For details, please refer to : <http://c.raqsoft.com/article/1571711202215>

➤ 1. esProc External library



After the model is created (or the model created by the intelligent modeling designer), the external library of intelligent modeling can be called through SPL for prediction. The SPL for Prediction:

	A	B
1	=ym_env()	/Initialize environment
2	=ym_load_pcf("titanic.pcf")	/Loading model file
3	=file("titanic_test.csv").import@cqt()	/Loading prediction data
4	=ym_predict(A2,A3)	/Execute prediction, return predicted result object
5	=ym_result(A4)	/Get predicted result sequence table
6	=ym_json(A4)	/When the prediction data is no less than 20 pieces, the model performance JSON information will be exported according to the prediction data evaluation.
7	>ym_close(A1)	/Close

A5

PassengerId	Survived	Pclass	Name	Sex	...
624	0	3	Hansen,...	male	...
625	0	3	Bowen, ...	male	...
626	0	1	Sutton, ...	male	...
627	0	2	Kirkland...	male	...
...

A6

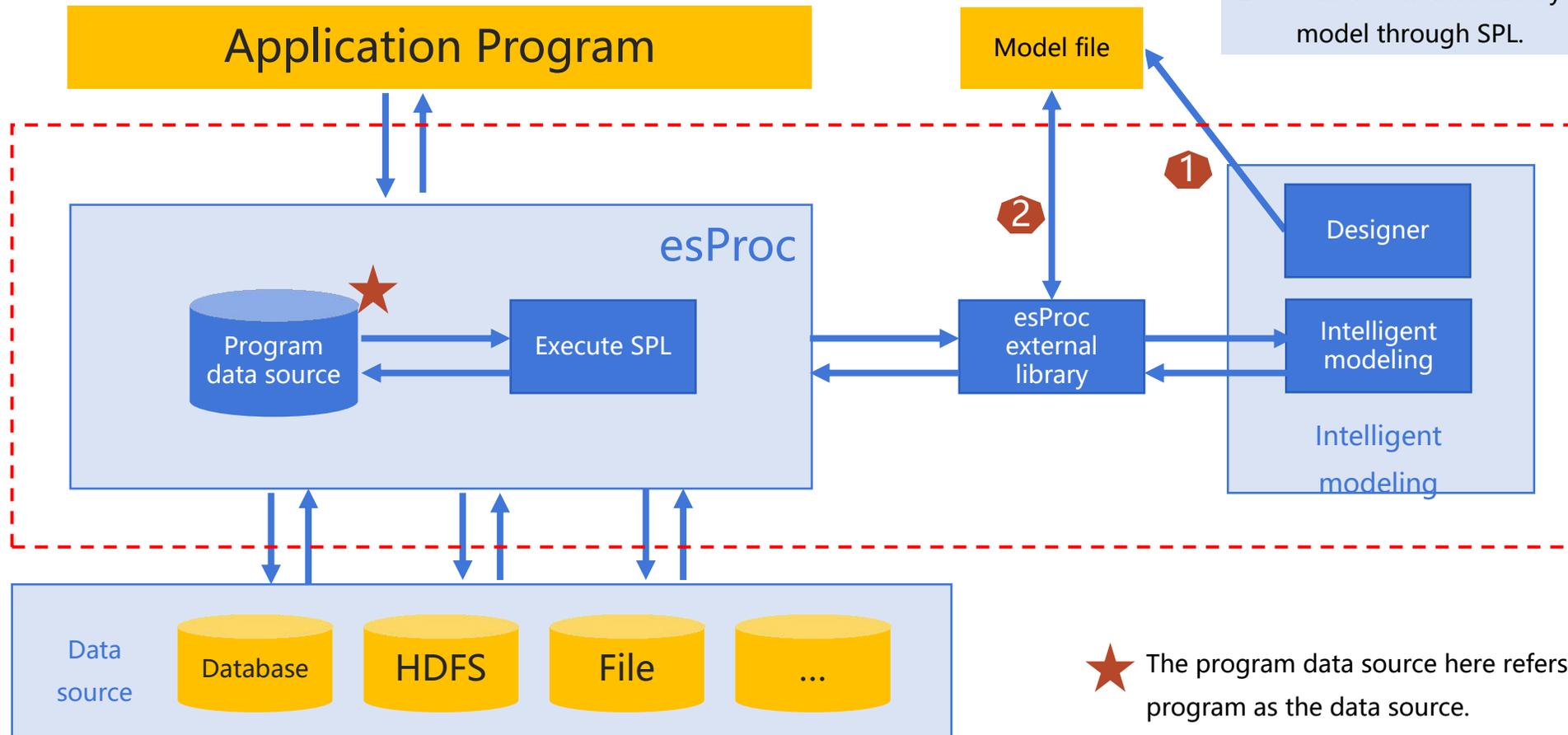
Value
{“Model-Performance”:{“GINI”:0.8369670542635659,“AUC”:0.9184835271317829,“KS”:0.6867732558139534,“ROC-Data\“:[“1-specificity\“:\“0.0\“,“sensitivity\“:\“0.020833333333333332\“]“,\“1-...

2. Integration architecture



There are two ways to create a model:

1. Use the intelligent modeling designer to create model file
2. Call the external library of esProc to create model through SPL.



THANKS

Innovation makes progress



www.raqsoft.com.cn