



# Altair Data Science Contest

By Aishani singh





# What is RapidMiner?

**RapidMiner is a Data Science platform by Altair that can be used for various task like preparation, text mining ,predictive analytics and machine learning (ML)**



# Lung cancer Prediction

- we are going to predict the lung cancer and dataset is taken from kaggle -which is an Online platform for data science
- task is to predict the if the person has a cancer or not.
- Link of kaggle :  
<https://www.kaggle.com/datasets/mysarahmadbhat/lung-cancer>

# About the Analysis

- Here , we are using RapidMiner to build a predictive model which will predict
- the lung cancer on dataset
- Key features of RapidMiner is GUI.
- there are 2 method to predict the build model.
  - Automodel
  - Using process Panel
- Another feature is drag and drop feature ,which is more easier for the data analyst
- there are various method and algorithms and operators are available in RapidMiner studio.

# RapidMiner features

- Drag-and-Drop Interface: Simplifies the creation of data science workflows with an intuitive, code-free, visual design.
- Extensive Operator Library: Provides over 1500 operators for tasks like data preprocessing, machine learning, and model evaluation.
- Auto Model: Automates model selection, hyperparameter tuning, and evaluation to quickly find the best-performing models.
- Real-Time Scoring: Enables real-time predictions and data scoring, facilitating fast deployment in production environments.
- Integration with Big Data Platforms: Connects seamlessly with big data sources like Hadoop and Spark for large-scale data processing.

# About Dataset

- Name of the dataset survey lung cancer.csv
- It has 309 records with 16 attributes.

## Steps to cleansing the dataset:

- Remove Duplicates
- Handle Missing values

## Importing database:

- click the import data
- select the location and file to be used for analysis



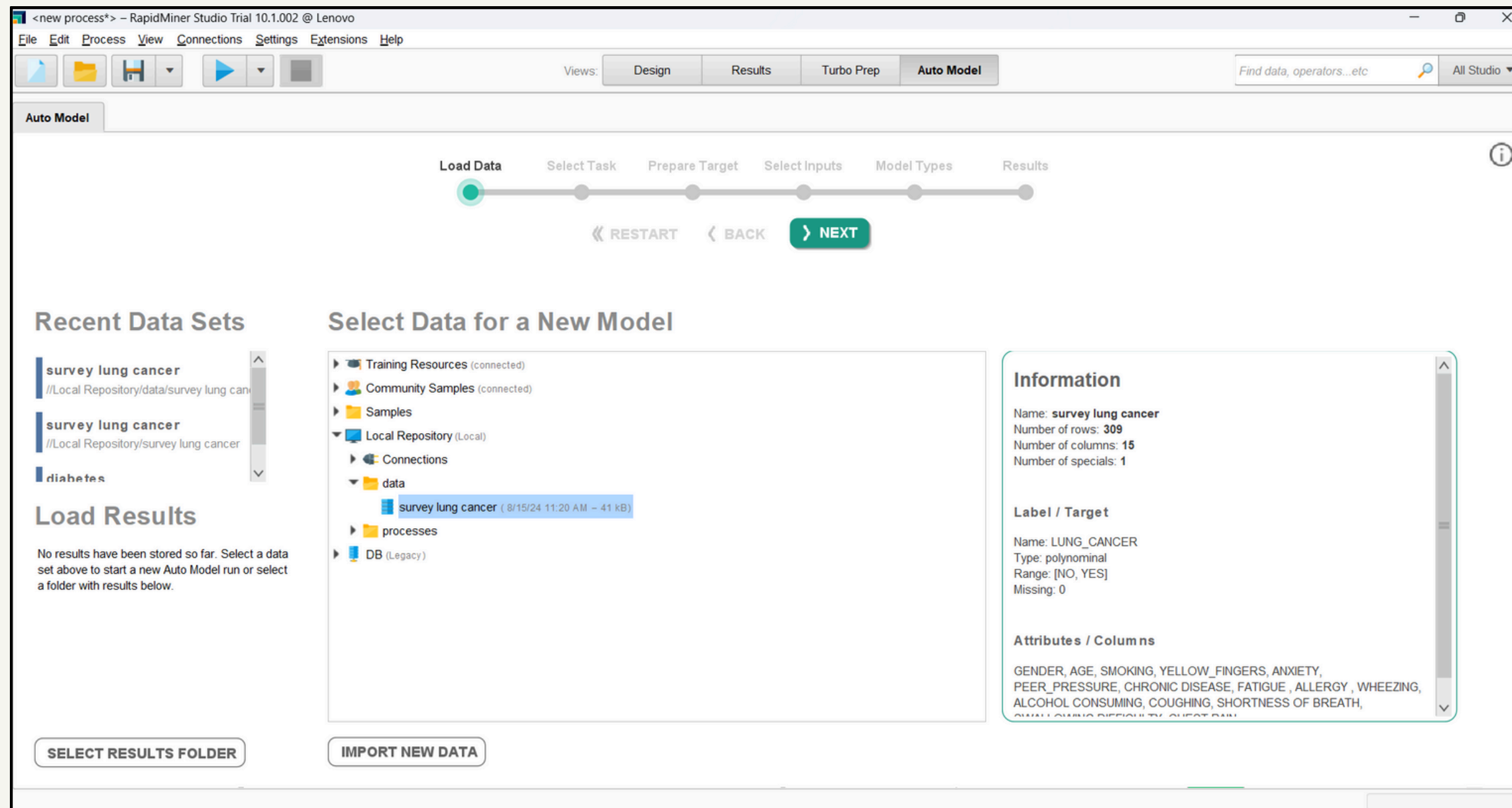
# Attributes

- GENDER
- LUNG\_CANCER
- CHEST PAIN
- SWALLOWING DIFFICULTY
- COUGHING
- ALCOHOL CONSUMING
- WHEEZING
- GENDER

- ALLERGY
- FATIGUE
- CHRONIC DISEASE
- PEER\_PRESSURE
- ANXIETY
- YELLOW\_FINGERS
- SMOKING
- AGE

# Predictive model using Automodel

## Step 1. Load Your Data





## Step 2. Transform the data

Remove attributes which are not contributed to the heart attack risk.

<new process\*> – RapidMiner Studio Trial 10.1.002 @ Lenovo

File Edit Process View Connections Settings Extensions Help

Views:

DesignResultsTurbo PrepAuto Model

Find data, operators...etc

All Studio ▾

Turbo Prep

Data Sets

+ LOAD DATA

survey lung cancer

//Local Repository/data/survey lung cancer

Rows: 309

Columns: 16

Last Change: None

survey lung cancer

Add new data sets on the left. Details for the selected data are shown below. You can change the data with the following actions.

TRANSFORM

CLEANSE

GENERATE

PIVOT

MERGE

MODELCHARTSCREATE PROCESSHISTORY

GENDER <small>Category</small>	AGE <small>Number</small>	SMOKING <small>Number</small>	YELLOW_FIN... <small>Number</small>	ANXIETY <small>Number</small>	PEER_PRES... <small>Number</small>	CHRONIC DI... <small>Number</small>	FATIGUE <small>Number</small>	ALLERGY <small>Number</small>	WHEEZING <small>Number</small>	ALCOHOL C... <small>Number</small>
M	69	1	2	2	1	1	2	1	2	2
M	74	2	1	1	1	2	2	2	1	1
F	59	1	1	1	2	1	2	1	2	1
M	63	2	2	2	1	1	1	1	1	2
F	63	1	2	1	1	1	1	1	2	1
F	75	1	2	1	1	2	2	2	2	1
M	52	2	1	1	1	1	2	1	2	2
F	51	2	2	2	2	1	2	2	1	1
F	68	2	1	2	1	1	2	1	1	1
M	53	2	2	2	2	2	1	2	1	2
F	61	2	2	2	2	2	2	1	2	1
M	72	1	1	1	1	2	2	2	2	2

309 rows - 16 columns (2 nominal, 14 numerical)

# Step 3. Load the Data and Select the task

<new process\*> - RapidMiner Studio Trial 10.1.002 @ Lenovo

File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep **Auto Model**

Find data, operators...etc All Studio ▾

**Auto Model**

Load Data **Select Task** Prepare Target Select Inputs Model Types Results

⏪ RESTART ⏪ BACK **> NEXT**

**Predict**  
Want to predict the values of a column?

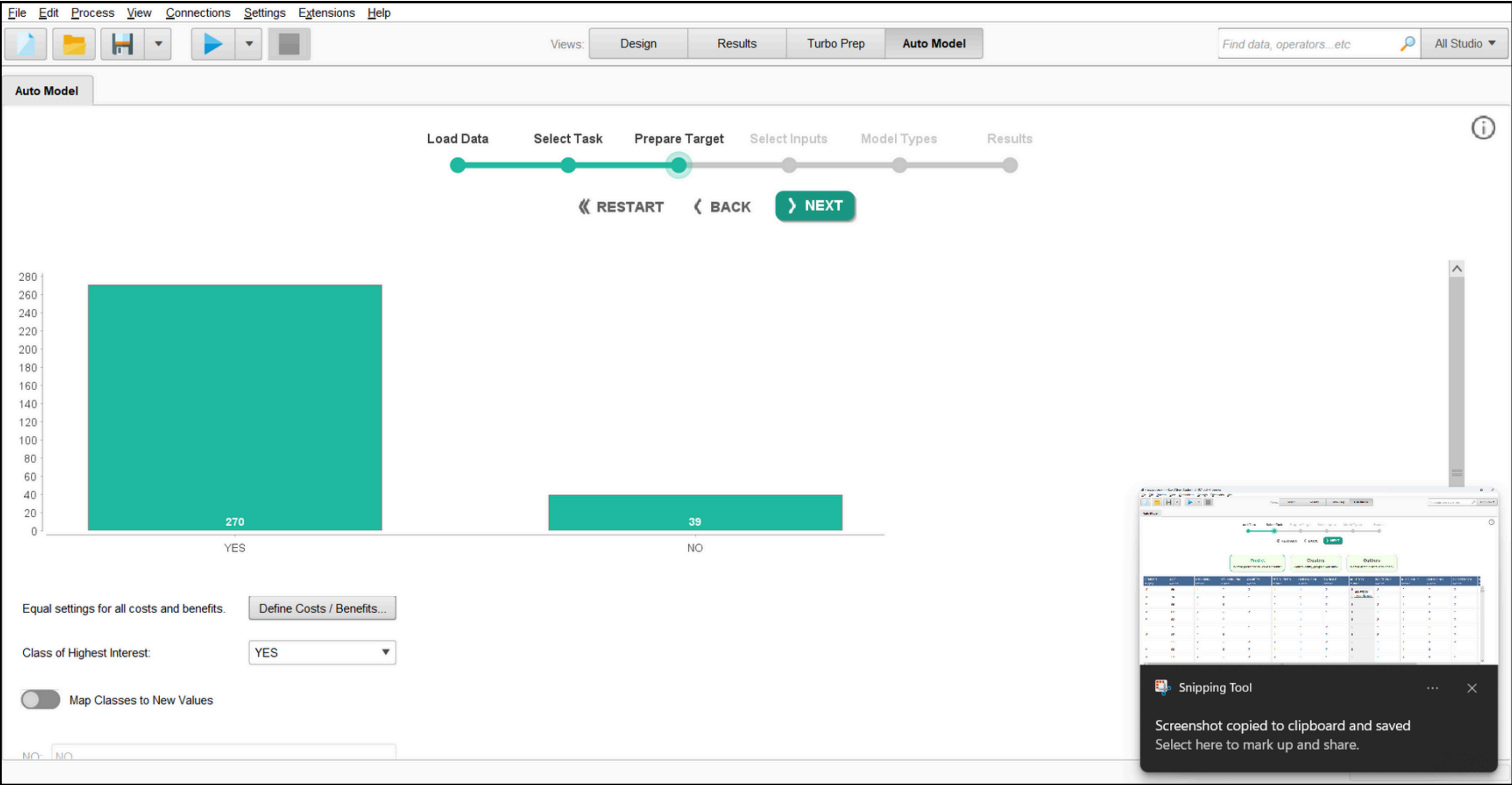
**Clusters**  
Want to identify groups in your data?

**Outliers**  
Want to detect outliers in your data?

GENDER Category	AGE Number	SMOKING Number	YELLOW_FIN... Number	ANXIETY Number	PEER_PRES... Number	CHRONIC DI... Number	FATIGUE Number	ALLERGY Number	WHEEZING Number	ALCOHOL C... Number	COUGHING Number	SHORTNESS ... Number	S... Number
M	69	1	2	2	1	1	2	1	2	2	2	2	
M	74	2	1	1	1	2	2	2	1	1	1	2	
F	59	1	1	1	2	1	2	1	2	1	2	2	
M	63	2	2	2	1	1	1	1	1	2	1	1	
F	63	1	2	1	1	1	1	1	2	1	2	2	
F	75	1	2	1	1	2	2	2	2	1	2	2	
M	52	2	1	1	1	1	2	1	2	2	2	2	
F	51	2	2	2	2	1	2	2	1	1	1	2	
F	68	2	1	2	1	1	2	1	1	1	1	1	
M	53	2	2	2	2	2	1	2	1	2	1	1	

309 rows - 16 columns (1 nominal, 14 numerical)

# Step 4.Prepare Target



# Step 5. Select Inputs

<new process\*> – RapidMiner Studio Trial 10.1.002 @ Lenovo

FileEditProcessViewConnectionsSettingsExtensionsHelp

Views: 

Design

Results

Turbo Prep

Auto Model

Find data, operators...etc

All Studio

Auto Model

Load Data

Select Task

Prepare Target

Select Inputs

Model Types

Results

RESTART

BACK

NEXT

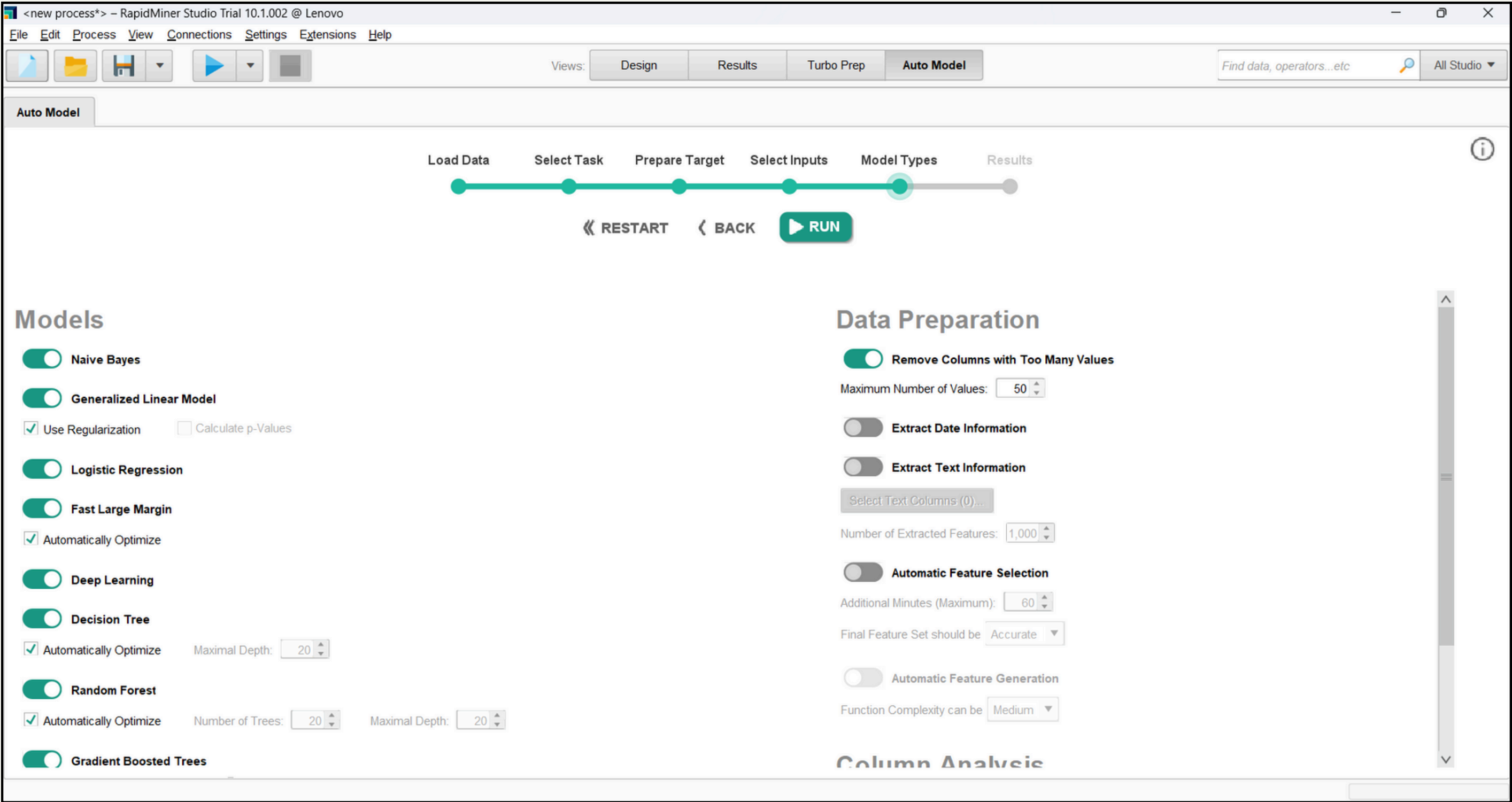
Selected: 15 / Total: 15

Select All

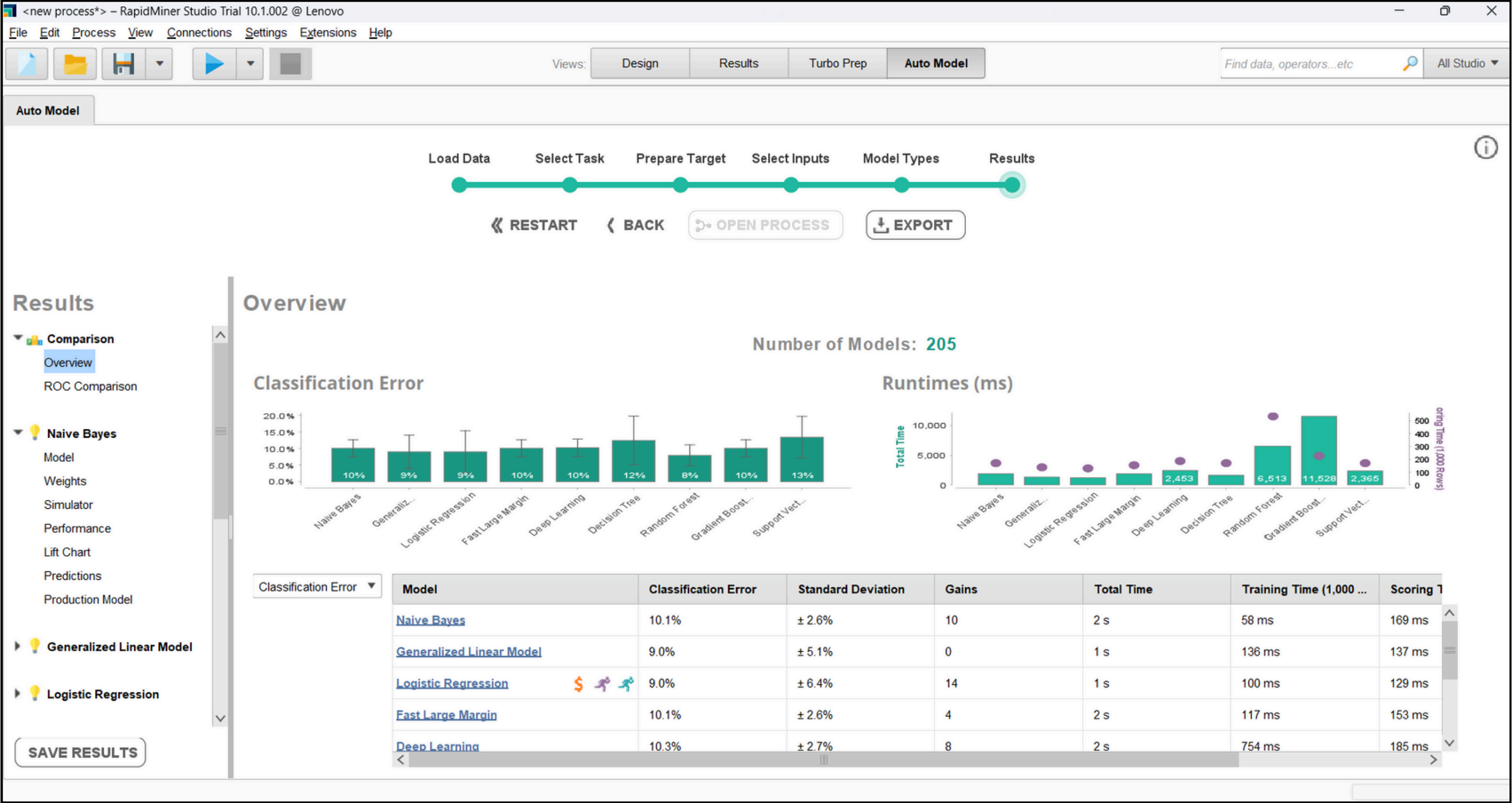
Deselect All

Selected	Status ↑	Quality	Name	Correlation	ID-ness	Stability	Missing	Text-ness
<input checked="" type="checkbox"/>	<div></div>	<div>C I S M T</div>	<div></div> GENDER	0.45%	0.65%	52.43%	0.00%	0.66%
<input checked="" type="checkbox"/>	<div></div>	<div>C I S M T</div>	<div></div> AGE	0.80%	12.62%	6.47%	0.00%	0.00%
<input checked="" type="checkbox"/>	<div></div>	<div>C I S M T</div>	SMOKING	0.34%	0.65%	56.31%	0.00%	0.00%
<input checked="" type="checkbox"/>	<div></div>	<div>C I S M T</div>	YELLOW_FINGERS	3.29%	0.65%	56.96%	0.00%	0.00%
<input checked="" type="checkbox"/>	<div></div>	<div>C I S M T</div>	ANXIETY	2.10%	0.65%	50.16%	0.00%	0.00%
<input type="checkbox"/>	<div></div>	<div>C I</div>						

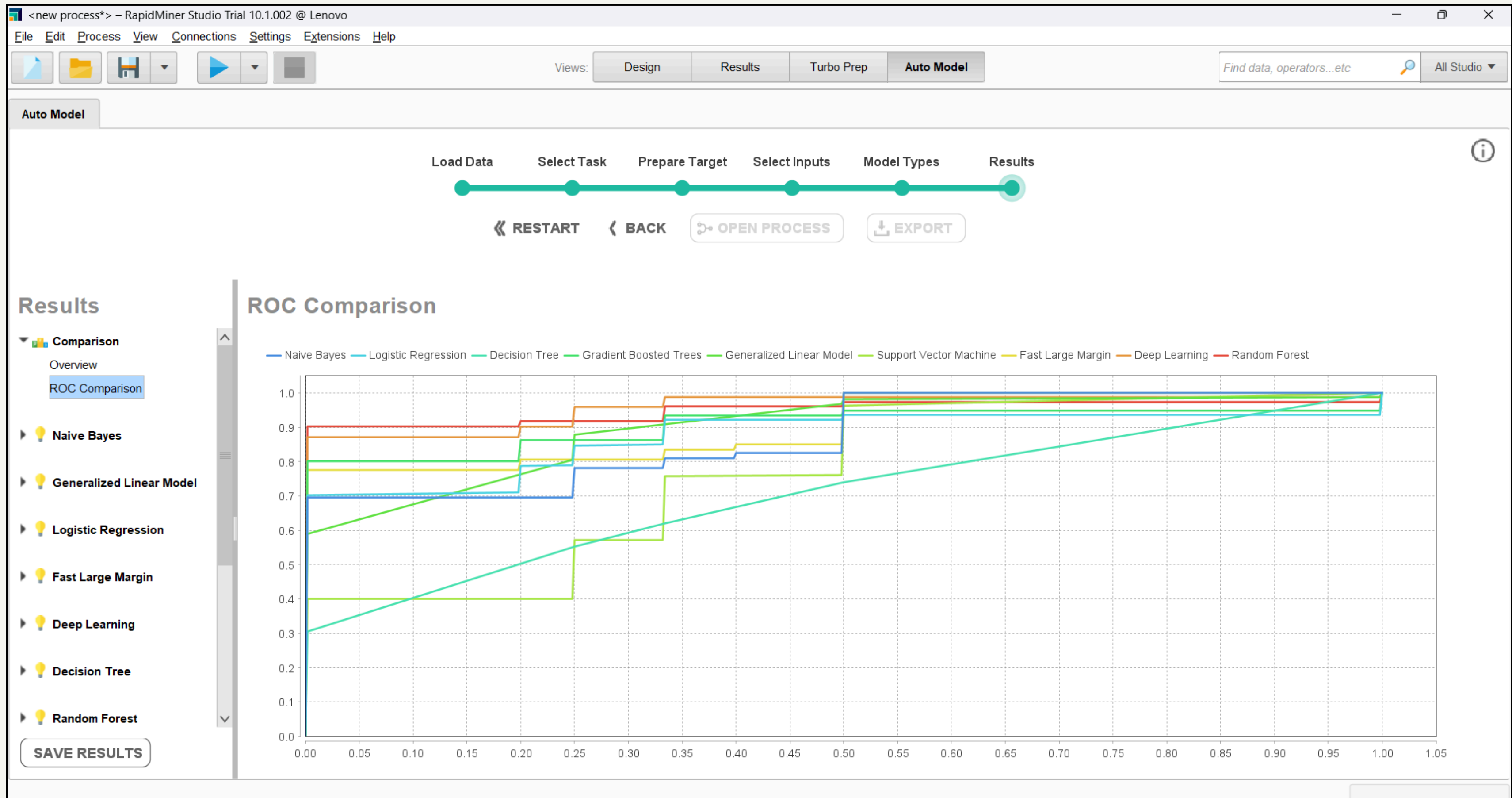
# Step 6. Select Model type



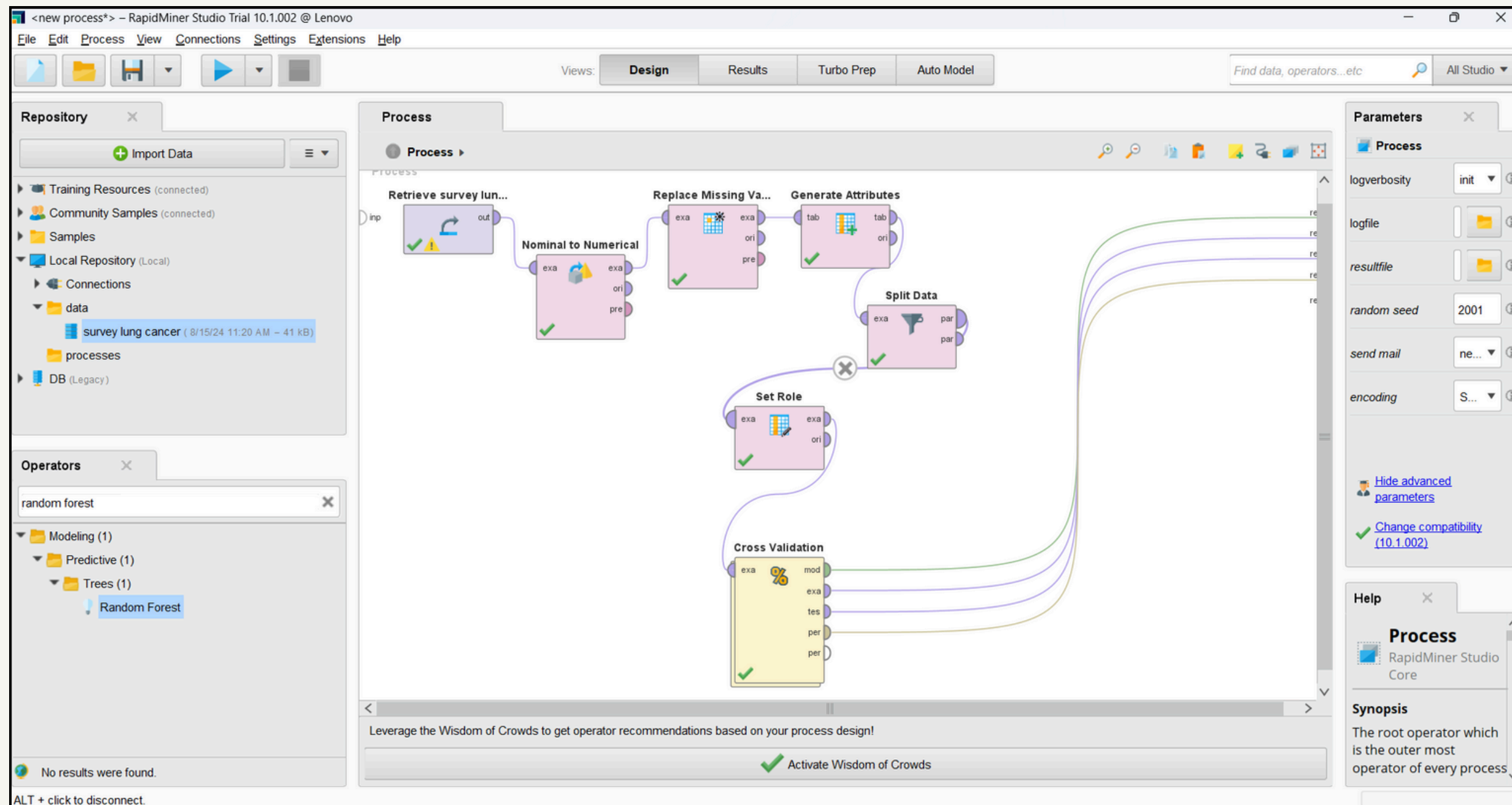
# Result Page





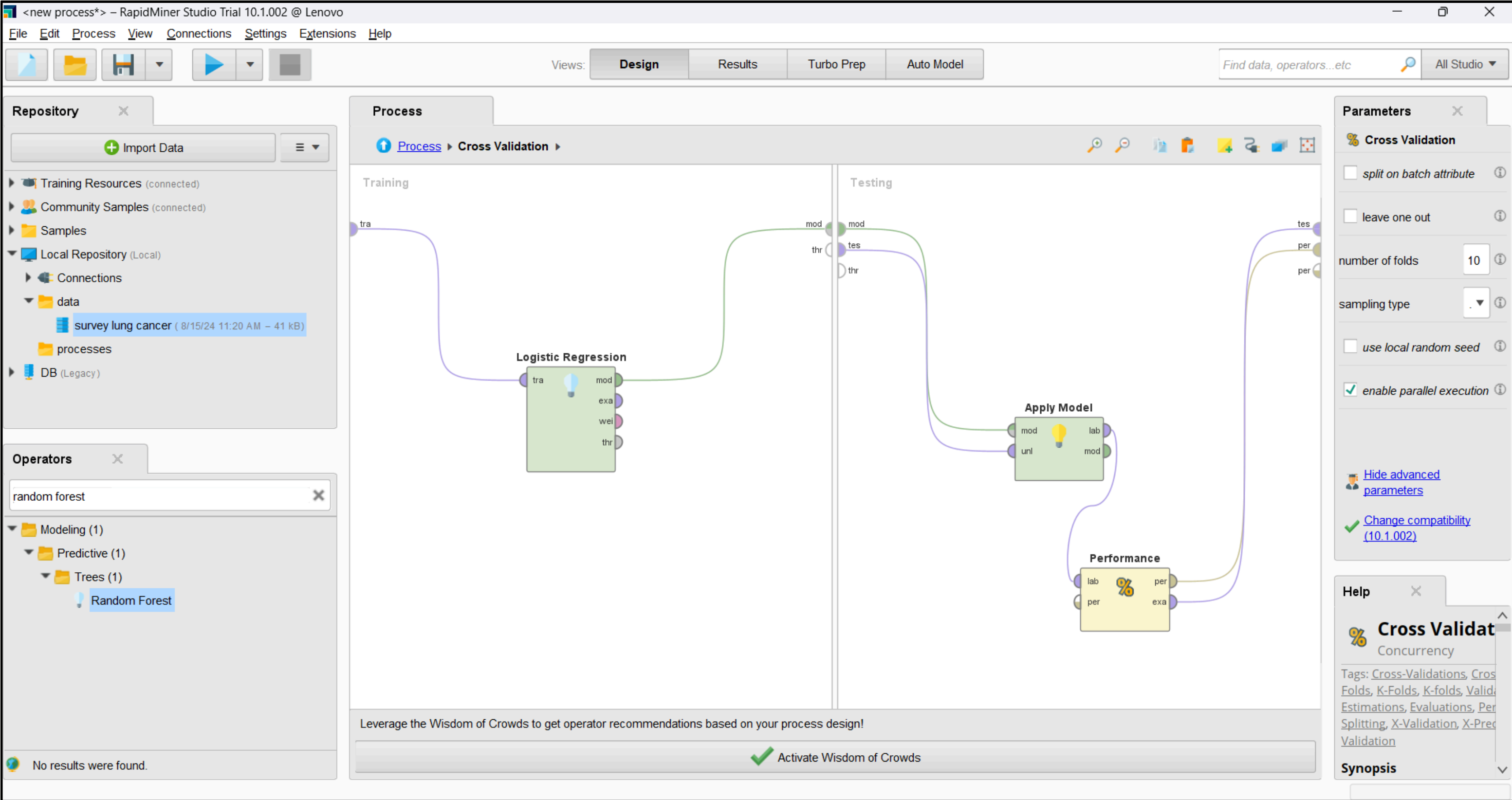


# Predictive model using Process Panel

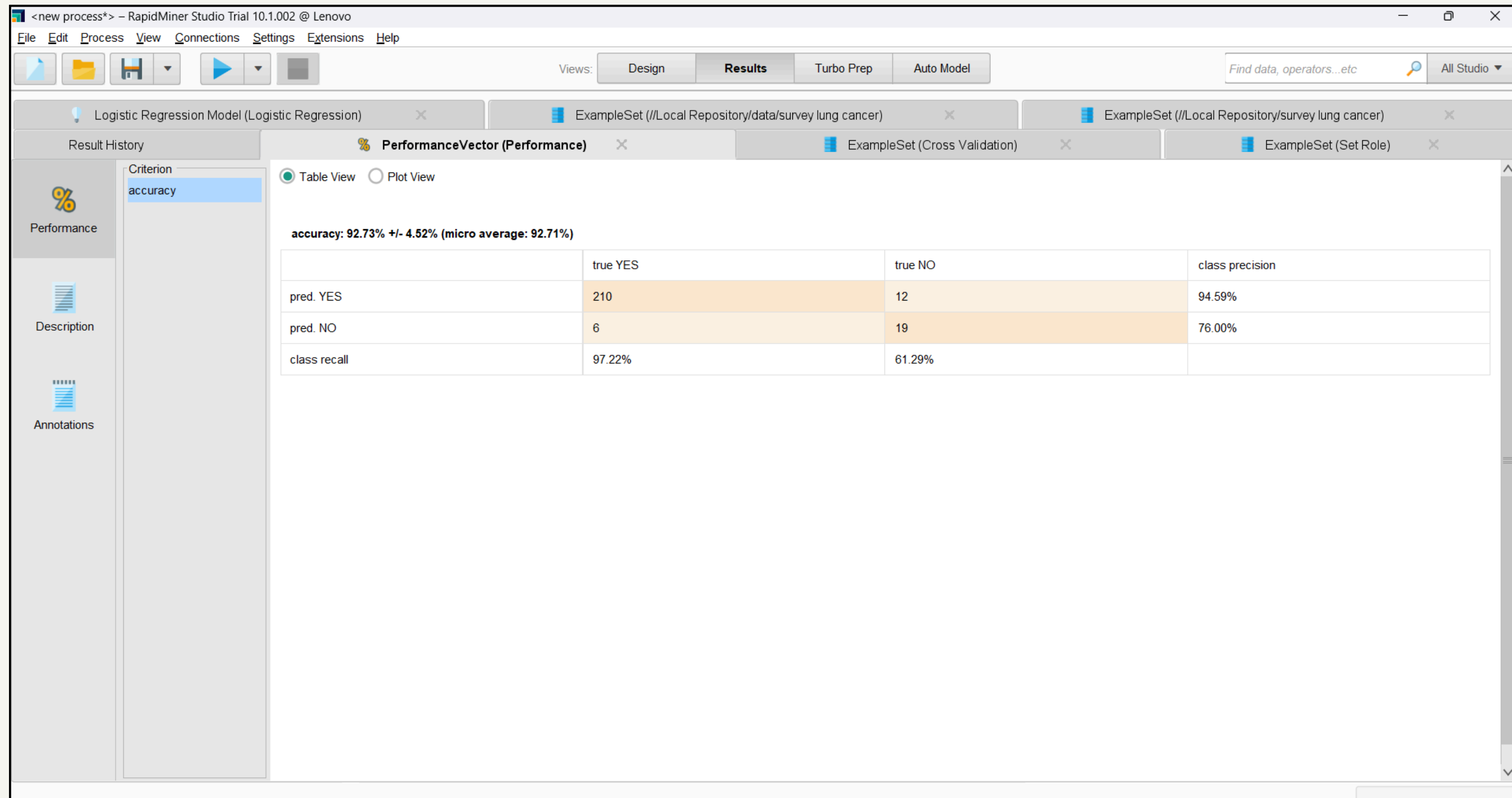




# In cross validation



# when we Run



<new process\*> – RapidMiner Studio Trial 10.1.002 @ Lenovo

FileEditProcessViewConnectionsSettingsExtensionsHelp

Views: 

Design

Results

Turbo Prep

Auto Model

Find data, operators...etc

All Studio

Logistic Regression Model (Logistic Regression)

ExampleSet (//Local Repository/data/survey lung cancer)

ExampleSet (//Local Repository/survey lung cancer)

Result History

PerformanceVector (Performance)

ExampleSet (Cross Validation)

ExampleSet (Set Role)

Data

Statistics

Visualizations

Annotations

Open in

Turbo Prep

Auto Model

Filter (247 / 247 examples):

all

Row No.	LUNG_CAN...	prediction(L...	confidence(...	confidence(...	GENDER = M	GENDER = F	AGE	SMOKING	YELLOW_FI...	ANXIETY	PEER_PRE...	CHRONIC D...	FATIGUE
1	YES	YES	0.040	0.960	1	0	52	2	1	1	1	1	2
2	YES	YES	0.008	0.992	1	0	58	2	1	1	1	1	2
3	YES	YES	0.008	0.992	1	0	60	2	1	1	1	1	2
4	YES	YES	0.001	0.999	0	1	53	2	2	2	1	2	1
5	YES	YES	0.066	0.934	1	0	62	2	1	2	1	1	1
6	NO	YES	0.195	0.805	0	1	56	1	1	1	1	2	1
7	YES	YES	0.077	0.923	1	0	52	2	1	1	2	1	2
8	YES	YES	0.000	1.000	1	0	62	2	2	1	1	2	1
9	YES	YES	0.000	1.000	0	1	70	1	1	2	2	2	2
10	YES	YES	0.345	0.655	1	0	58	2	1	1	1	1	1
11	YES	YES	0.042	0.958	0	1	49	1	1	1	2	2	1
12	NO	NO	0.634	0.366	0	1	63	1	1	1	1	2	2
13	YES	YES	0.000	1.000	0	1	47	2	2	1	2	2	2
14	YES	YES	0.010	0.990	1	0	67	2	1	2	1	1	2
15	YES	YES	0.000	1.000	1	0	67	1	2	2	2	1	2
16	YES	YES	0.378	0.622	1	0	54	2	1	1	1	1	1

ExampleSet (247 examples, 4 special attributes, 18 regular attributes)