

Правительство Российской Федерации  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ  
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»  
(НИУ ВШЭ)

Московский институт электроники и математики им. А.Н. Тихонова

ПРАКТИЧЕСКАЯ РАБОТА № 3

ТЕМА РАБОТЫ  
«Визуализации в RapidMiner»

Москва, 2025

## Практическая работа 3: Визуализации в RapidMiner

Практическая работа 3: Визуализации в RapidMiner.....	2
Цель работы.....	2
Целевая аудитория.....	3
Идея и концепция.....	3
Содержание практической работы.....	3
Введение в RapidMiner.....	3
О наборе данных и задаче работы.....	4
Работа с данными.....	4
Загрузка набора данных.....	4
Использование TurboPrep для подготовки данных.....	5
Использование Transform, Cleanse, Generate, Pivot, Merge.....	8
Визуализация данных.....	8
Примеры визуализаций, которые необходимо создать.....	8
1. Гистограмма (Histogram) для Age.....	8
2. Столбчатая диаграмма (Bar Chart) для Pclass или Sex.....	8
3. Круговая диаграмма (Pie Chart) для Pclass.....	9
4. Диаграмма рассеяния (Scatter Plot) для Age и Fare, цвет по Survived.....	10
5. Box Plot (ящик с усами) для Fare по Embarked.....	10
6. Тепловая карта (Heatmap) корреляций для числовых признаков.....	11
Настройки графиков.....	11
Автоматизированное создание модели с помощью AutoModel.....	12
Анализ результатов.....	12
Приобретенные навыки.....	14
Обобщенная задача для выполнения индивидуального варианта.....	14
Распределение вариантов.....	16

### Цель работы

Познакомиться с расширенными возможностями визуализации данных в RapidMiner. Студенты научатся создавать различные типы диаграмм (гистограммы, круговые, диаграммы рассеяния, пузырьковые, box-plot), настраивать их внешний вид, а также использовать визуализацию для анализа результатов моделей. Главная задача – понять, как визуальный

анализ помогает глубже понять структуру данных и качество построенных моделей машинного обучения.

## Целевая аудитория

Студенты курса по информатике и анализу данных, имеющие базовые знания о работе с RapidMiner и предварительной обработке данных. Эта работа подходит для тех, кто уже знаком с основами загрузки и очистки данных (из предыдущих работ) и стремится освоить методы визуального анализа результатов.

## Идея и концепция

Для демонстрации возможностей визуализации используется набор данных «Titanic», содержащий информацию о пассажирах «Титаника». Этот набор данных широко используется благодаря сочетанию числовых и категориальных признаков. В нем представлены такие столбцы, как возраст пассажира, пол, класс билета, стоимость билета, порт посадки, а также информация о том, выжил ли пассажир.

Используя данный набор данных, студенты смогут:

- Визуально исследовать распределения признаков: возраст, стоимость билета (Fare), классы обслуживания.
- Изучать соотношение категориальных признаков (пол, класс, порт посадки) с целевой переменной (выжил/не выжил).
- Оценить корреляции между численными признаками и понять, какие факторы, возможно, повлияли на выживаемость.
- После предварительного анализа данных будет продемонстрирован процесс использования AutoModel для построения модели классификации (предсказания выживаемости) с последующей визуальной оценкой её результатов.

## Содержание практической работы

### Введение в RapidMiner

Краткое напоминание, что RapidMiner – это инструмент для анализа данных, включающий средства предобработки, визуализации и автоматизированного построения моделей.

- При необходимости установить RapidMiner Studio.
- Ознакомиться с панелью инструментов и вкладками «Design», «Results», Turbo Prep», «Auto Model».
- Убедиться в наличии файла с данными Titanic (csv-файл, доступный, например, [на Kaggle](#) или в локальном хранилище).

О наборе данных и задаче работы

Датасет «Titanic» содержит следующие столбцы (атрибуты):

- Survived: бинарный признак (0 – не выжил, 1 – выжил).
- Pclass: класс пассажира (1 – первый класс, 2 – второй класс, 3 – третий класс).
- Name: имя пассажира (текстовый признак, можно не использовать для визуализаций).
- Sex: пол пассажира (male/female).
- Age: возраст (числовой признак).
- SibSp: количество братьев/сестер/супругов на борту.
- Parch: количество родителей/детей на борту.
- Ticket: номер билета (текстовый, можно не использовать для визуализаций).
- Fare: стоимость билета (числовой признак).
- Cabin: номер каюты (может содержать много пропусков).
- Embarked: порт посадки (C – Cherbourg, Q – Queenstown, S – Southampton).

Задача: с помощью визуализации понять структуру данных Titanic, увидеть различия между группами пассажиров и их влияние на выживаемость. После этого построить модель машинного обучения (классификацию), оценить ее качество и интерпретировать результаты визуально.

Работа с данными

Загрузка набора данных

- Откройте RapidMiner Studio.
- Нажмите "Create New Process".
- Перетащите оператор "**Read CSV**" на рабочее пространство.
- Загрузите файл titanic.csv (предварительно скачанный).

- Подключите выход **Read CSV** к **Result** и нажмите **Run**, чтобы просмотреть данные.

Result History ExampleSet (//Local Repository/Titanic-Dataset) X

Open in Turbo Prep Auto Model Interactive Analysis Filter (891 / 891 examples): all

Row No.	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	1	0	3	Braund, Mr. ...	male	22	1	0	A/5 21171	7.250	?	S
2	2	1	1	Cummings, Mr...	female	38	1	0	PC 17599	71.283	C85	C
3	3	1	3	Heikkinen, ...	female	26	0	0	STON/O2. 3...	7.925	?	S
4	4	1	1	Futrelle, Mrs...	female	35	1	0	113803	53.100	C123	S
5	5	0	3	Allen, Mr. Wi...	male	35	0	0	373450	8.050	?	S
6	6	0	3	Moran, Mr. J...	male	?	0	0	330877	8.458	?	Q
7	7	0	1	McCarthy, M...	male	54	0	0	17463	51.862	E46	S
8	8	0	3	Palsson, Ma...	male	2	3	1	349909	21.075	?	S
9	9	1	3	Johnson, Mr...	female	27	0	2	347742	11.133	?	S
10	10	1	2	Nasser, Mrs...	female	14	1	0	237736	30.071	?	C
11	11	1	3	Sandstrom, ...	female	4	1	1	PP 9549	16.700	G6	S
12	12	1	1	Bonnell, Mis...	female	58	0	0	113783	26.550	C103	S
13	13	0	3	Saunders, ...	male	20	0	0	A/5. 2151	8.050	?	S
14	14	0	3	Andersson, ...	male	39	1	5	347082	31.275	?	S
15	15	0	3	Vestrom, Mi...	female	14	0	0	350406	7.854	?	S
16	16	1	2	Hewlett, Mrs...	female	55	0	0	248706	16	?	S

ExampleSet (891 examples, 0 special attributes, 12 regular attributes)

рис. 1: Просмотр датасета после загрузки

Использование TurboPrep для подготовки данных

- Нажмите кнопку **"Turbo Prep"** в верхней панели инструментов.
- Загрузите набор данных Titanic.
- Проверьте данные на пропуски (особенно в столбцах Age, Cabin).
- При необходимости используйте функции **Cleanse** для замены пропусков в Age на среднее значение.

Name | Type | Missing | Statistics | Filter (12 / 12 attributes): Search for Attributes

PassengerId	Integer	0	Min 1, Max 891, Average 446
Survived	Integer	0	Min 0, Max 1, Average 0.384
Pclass	Integer	0	Min 1, Max 3, Average 2.309
Name	Polynomial	0	Least van Melk [...], Most Abbing, Mr. Anthony (1), Values Abbing, Mr. Anthony (1), Abbott, [...] re Edward (1), ...[889]
Sex	Polynomial	0	Least female (314), Most male (577), Values male (577), female (314)
Age	Real	177	Min 0.420, Max 80, Average 29.699
SibSp	Integer	0	Min 0, Max 8, Average 0.523
Parch	Integer	0	Min 0, Max 6, Average 0.382

рис. 2: Просмотр информации о датасете, включающей ключевые статистические показатели

### Cleanse

1 column selected

AUTO CLEANSING

REMOVE LOW QUALITY

REMOVE CORRELATED

REPLACE MISSING

NORMALIZATION

DISCRETIZATION

DUMMY ENCODING

PCA

REMOVE DUPLICATES

### Titanic-Dataset

Select a column to clean (hold Shift for selecting a range of columns; Ctrl for (de-)selecting multiple columns; Alt to select all columns of the same type; Ctrl+A for all columns). ...

COMMIT CLEANSE CANCEL

UNDO SHOW HISTORY

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
Number	Number	Number	Category	Category	Number	Number	Number	Category	Number
1	0	3	Braund, Mr. O...	male	22	1	0	A/5 21171	7.250
2	1	1	Cumings, Mrs....	female	38	1	0	PC 17599	71.283
3	1	3	Heikkinen, Mi...	female	26	0	0	STON/O2. 31...	7.925
4	1	1	Futrelle, Mrs. ...	female	35	1	0	113803	53.100
5	0	3	Allen, Mr. Willi...	male	35	0	0	373450	8.050
6	0	3	Moran, Mr. Ja...	male	?	0	0	330877	8.458
7	0	1	McCarthy, Mr....	male	54	0	0	17463	51.862
8	0	3	Palsson, Mast...	male	2	3	1	349909	21.075
9	1	3	Johnson, Mrs. ...	female	27	0	2	347742	11.133
10	1	2	Nasser, Mrs. ...	female	14	1	0	237736	30.071
11	1	3	Sandstrom, Mi...	female	4	1	1	PP 9549	16.700

891 rows - 12 columns (5 nominal, 7 numerical)

рис. 3: Заполнение пропусков в столбцах с помощью "Replace Missing"

- Вы увидите, что признак «Пол» (Sex) не является числовым типом данных, поэтому мы не сможем использовать этот признак в модели. Для решения этой проблемы стоит использовать Replace. Применим бинарное кодирование, заменив female (женский пол) на 1, а male (мужской пол) на 0.

### Transform

1 column selected

REMOVE

COPY

FILTER

RANGE

SAMPLE

SORT

REPLACE

male

0

Use regular expressions

APPLY

### Titanic-Dataset

Select columns to transform (hold Shift for selecting a range of columns; Ctrl for (de-)selecting multiple columns; Alt to select all columns of the same type; Ctrl+A for all columns). ...

COMMIT TRANSFORMATION CANCEL

PassengerId	Survived	Pclass	Name	Sex	Age
Number	Number	Number	Category	Category	Number
1	0	3	Braund, Mr. O...	male	22
2	1	1	Cumings, Mrs....	1	38
3	1	3	Heikkinen, Mi...	1	26
4	1	1	Futrelle, Mrs. ...	1	35
5	0	3	Allen, Mr. Willi...	male	35
6	0	3	Moran, Mr. Ja...	male	29.699
7	0	1	McCarthy, Mr....	male	54
8	0	3	Palsson, Mast...	male	2
9	1	3	Johnson, Mrs. ...	1	27
10	1	2	Nasser, Mrs. ...	1	14
11	1	3	Sandstrom, Mi...	1	4

891 rows - 12 columns (5 nominal, 7 numerical)

рис. 4: Использование оператора Replace для бинарного кодирования категориальных столбцов

- Обозначьте целевую переменную. В данном случае, мы предсказываем выживаемость, поэтому таргетом будет являться столбец «Survived». Делаем это с помощью оператора «Set role».

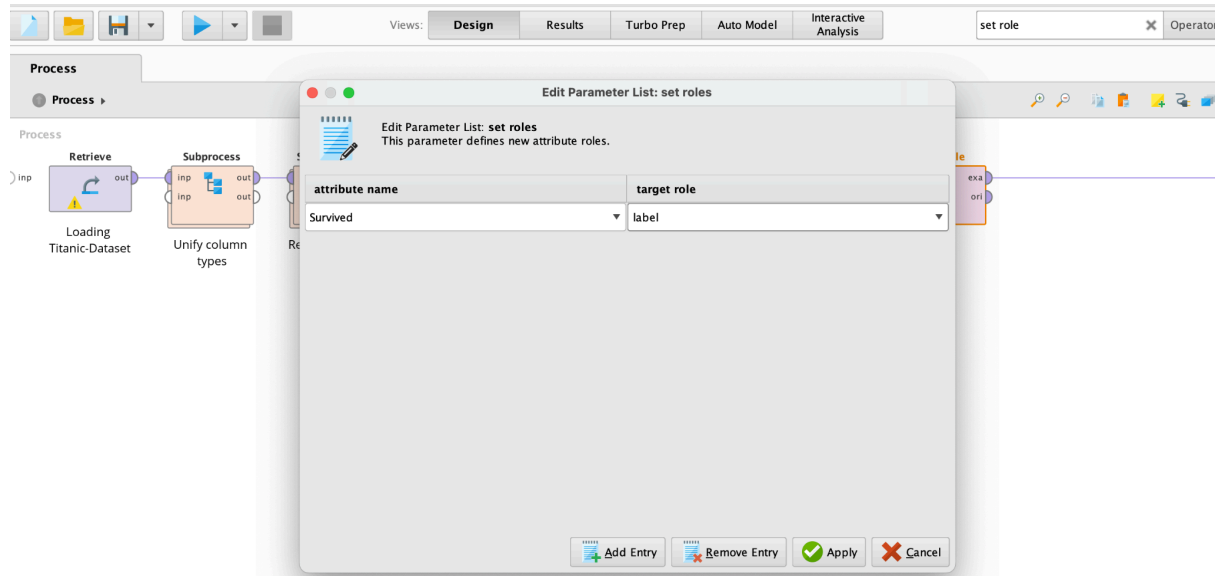


рис. 5: Выбор целевой переменной с помощью оператора "Set role"

- Нормализуйте числовые признаки для облегчения сравнения на графиках.

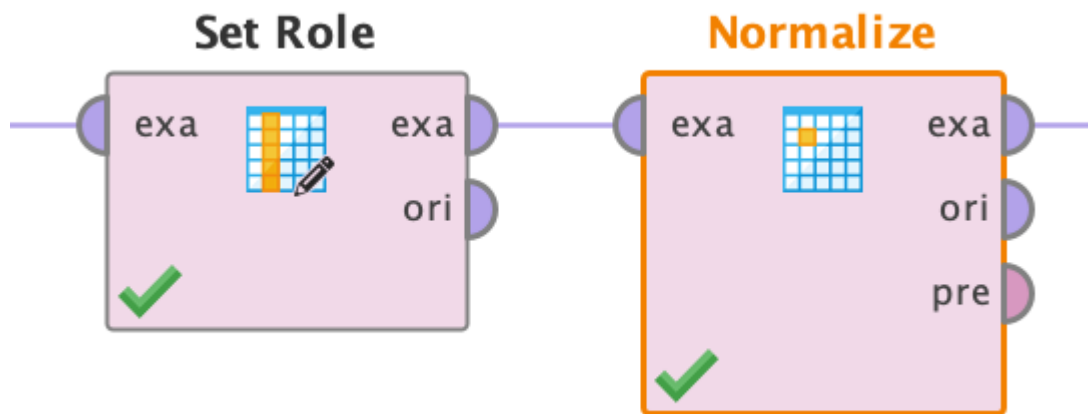


рис. 6: Добавление оператора "Normalize" в цепочку операторов для масштабирования значений и приведения к единой шкале

- Нажмите "Commit Cleanse", чтобы применить изменения.

## Использование Transform, Cleanse, Generate, Pivot, Merge

- **Transform:** при желании можно создать новый признак, например, категоризировать возраст по группам (молодые, взрослые, пожилые).
- **Cleanse:** уже применено для пропусков и нормализации.
- **Generate:** мсоздать новый признак, например, является ли возраст детским ( $< 16$  лет).
- **Pivot, Merge:** для данного датасета не обязательны.

## Визуализация данных

Откройте вкладку с подготовленными данными (TurboPrep → Charts).

Примеры визуализаций, которые необходимо создать

1. *Гистограмма (Histogram) для Age*
  - Plot type: Histogram
  - Value column: Age
  - Настройка: измените количество корзин (bins) для более детального анализа распределения возрастов.
  - Это позволит оценить распределение пассажиров по возрастным группам.
2. *Столбчатая диаграмма (Bar Chart) для Pclass или Sex*
  - Plot type: Bar
  - Value column: Pclass
  - X-axis column: Survived
  - Аналогично построить столбчатый график для Sex.
  - Это покажет, какова зависимость выживаемости от типа каюты или пола.

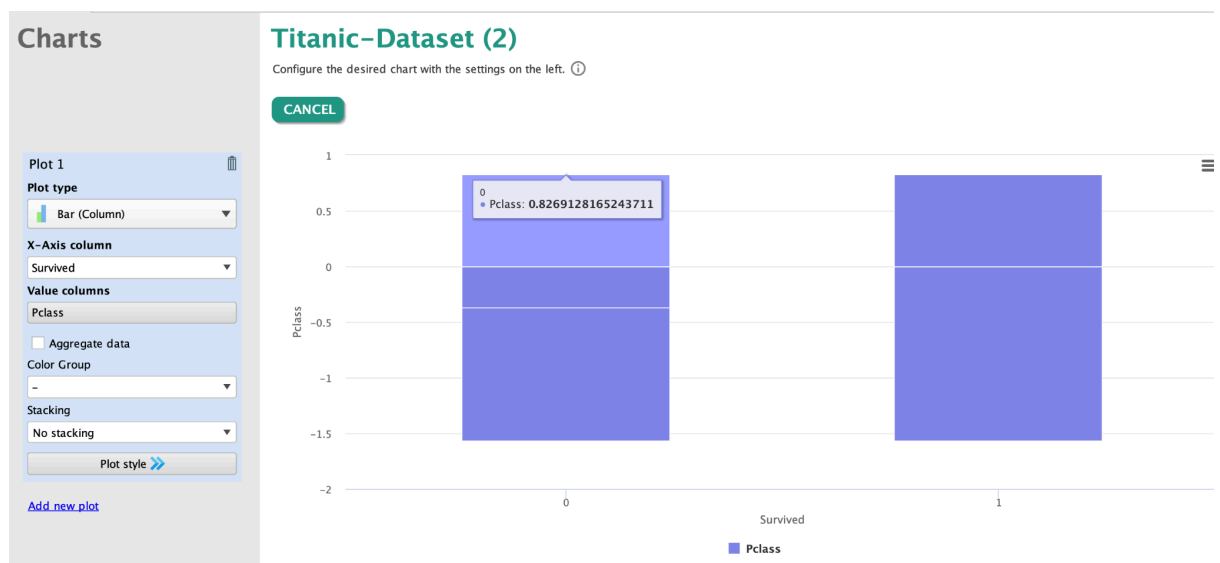




рис. 7: Построение столбчатой диаграммы "Зависимость выживаемости от типы каюта"

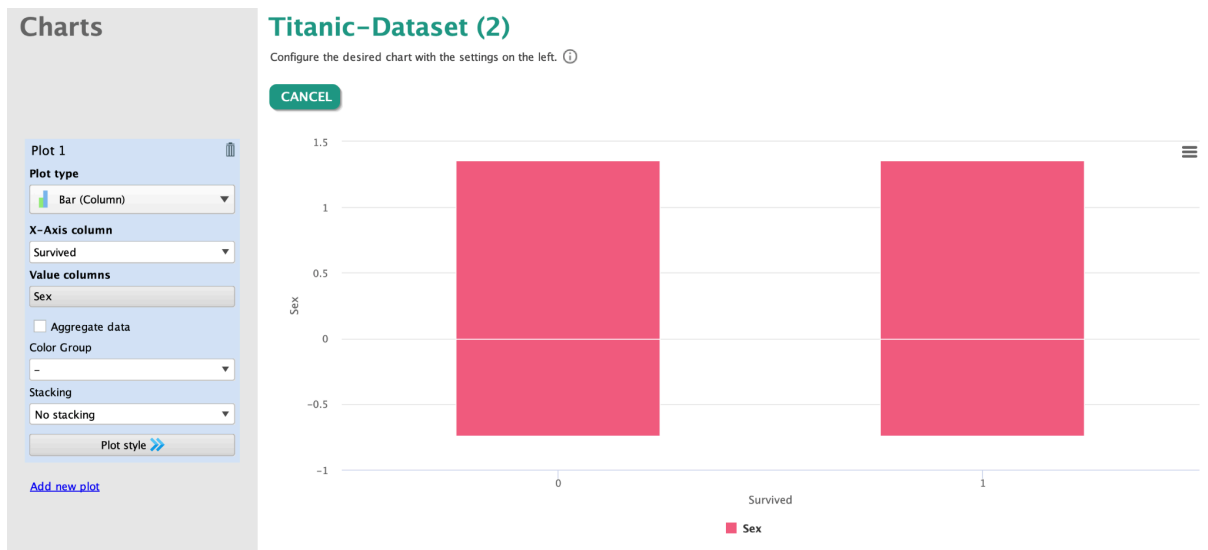


рис. 8: Построение столбчатой диаграммы "Зависимость выживаемости от типы от пола"

### 3. Круговая диаграмма (Pie Chart) для Pclass

- Plot type: Pie/Donut
- Value column: Survived
- Aggregation: Count
- Group by: Pclass
- Позволяет оценить долю выживших пассажиров в разных классах обслуживания (1-й, 2-й, 3-й классы). Цветовую схему можно изменить для лучшей наглядности. На графике будет видно, что выживаемость премиального класса выше.

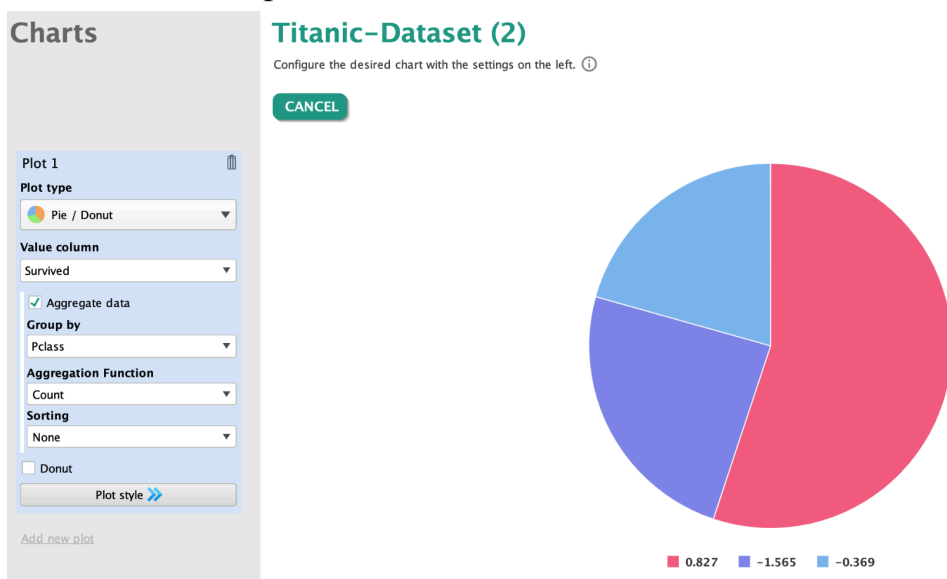


рис. 9: Построение круговой диаграммы "Доля выживших пассажиров в разных"

классах обслуживания"

4. Диаграмма рассеяния (Scatter Plot) для Age и Fare, цвет по Survived

- Plot type: Scatter/Bubble
- X-Axis column: Age
- Value columns: Fare
- Color by: Survived
- Определите, помогают ли эти признаки различать выживших и погибших. Измените форму и размер точек для лучшего восприятия.

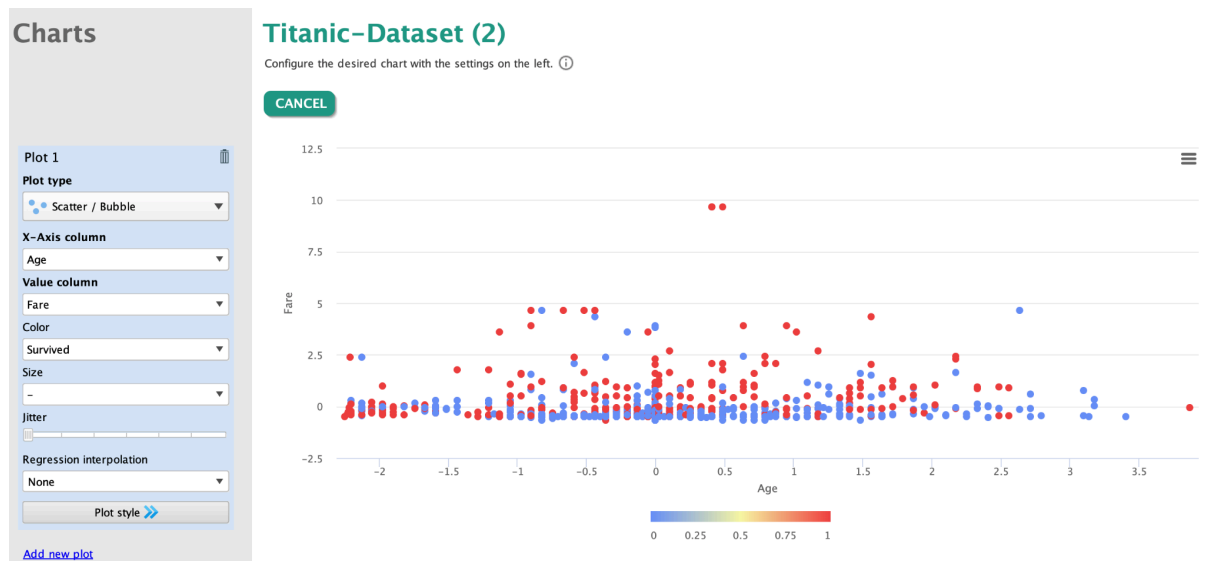


рис. 10: Построение диаграммы рассеяния

5. Box Plot (ящик с усами) для Fare по Embarked

- Plot type: Box
- Value column: Fare
- Group by: Embarked (порт отправления)
- Показывает, как стоимость билета распределена по портам отправления, есть ли сильные выбросы.

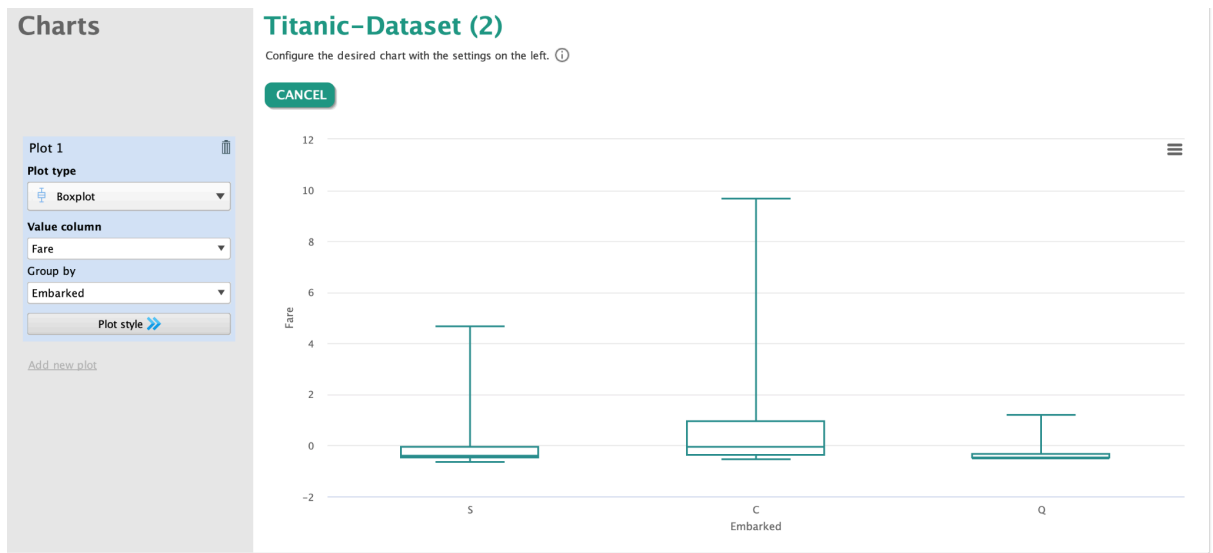


рис. 11: Построение ящика с усами (Box Plot)

#### 6. Тепловая карта (Heatmap) корреляций для числовых признаков

- Plot type: Heatmap
- Выберите матрицу корреляций между Age, Fare, SibSp, Parch.
- Показывает степень взаимосвязи признаков друг с другом, что важно при выборе признаков для модели.

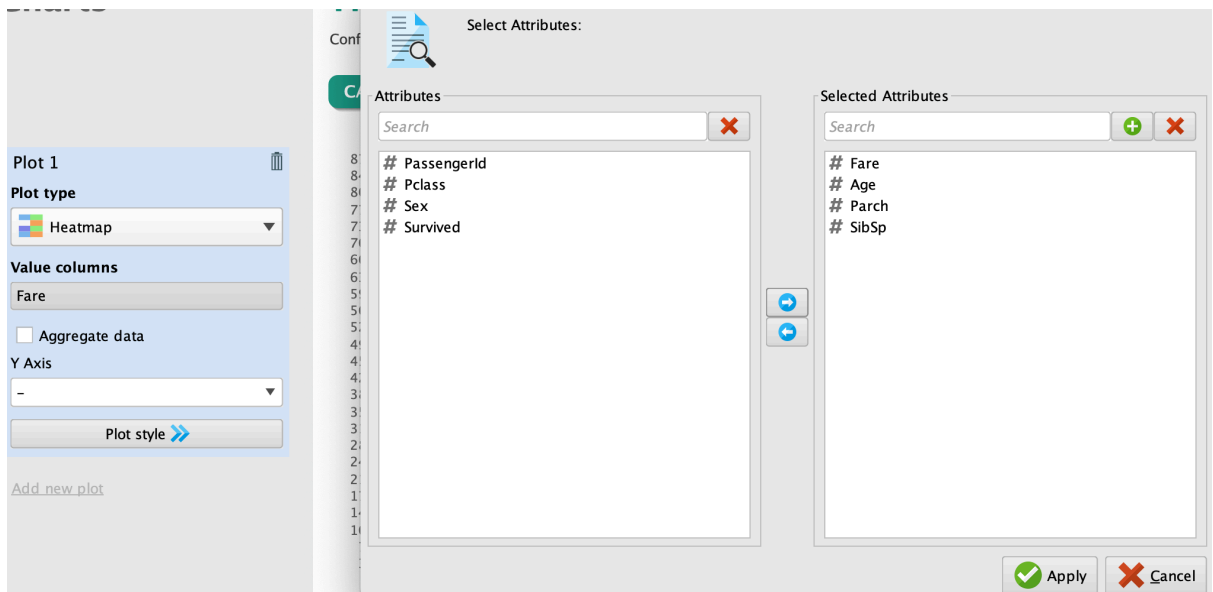


рис. 12: Построение тепловой карты: выбор признаков

### Настройки графиков

- Изменяйте заголовки осей (Axis Title), подписи к легенде (Legend), цветовые схемы (Color Settings).
- Экспериментируйте с фильтрами: например, отображать только пассажиров определенного класса или пола.

- Сортируйте столбчатые диаграммы по убыванию или возрастанию значений.

### Автоматизированное создание модели с помощью AutoModel

- Нажмите **"AutoModel"** в верхней панели.
- Загрузите подготовленный набор данных Titanic.
- Выберите задачу классификации с целевым признаком Survived.

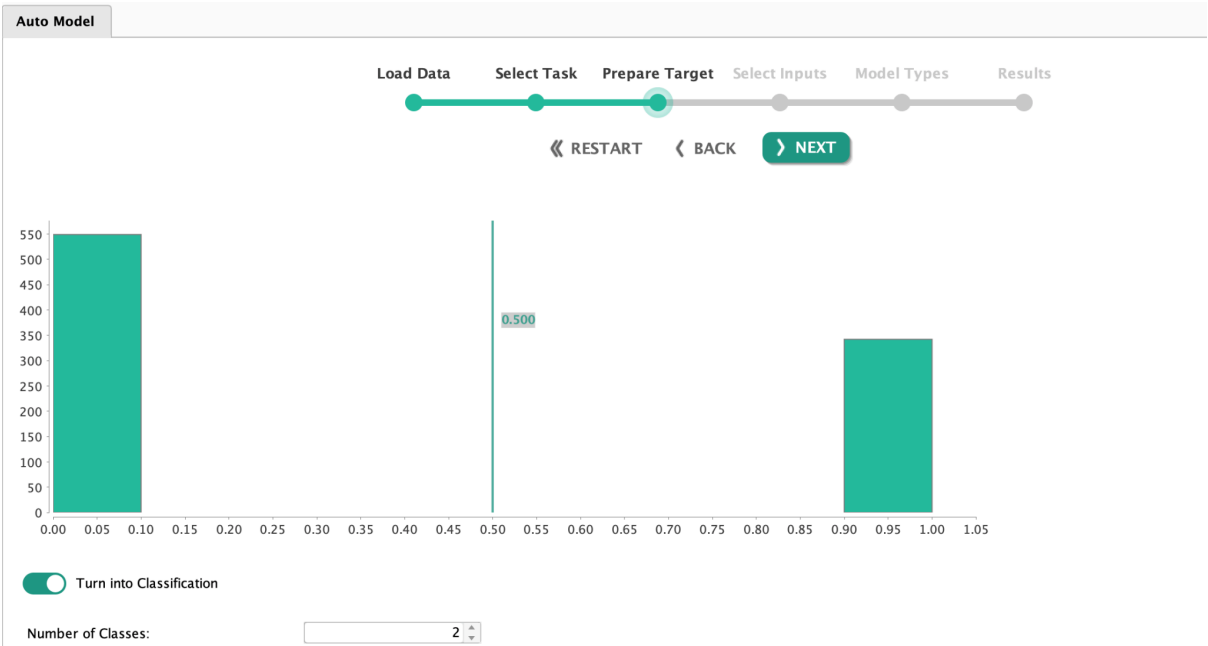


рис. 13: Использование AutoModel

- Выберите модели: Decision Tree, Logistic Regression, Random Forest, Naive Bayes.
- Нажмите **"Run"**, чтобы построить модели.

### Анализ результатов

- Перейдите во вкладку Results после работы AutoModel.
- Сравните точность (accuracy).

Model Category	Classification Error Number	Standard Deviation Number	Gains Number	Total Time Number	Training Time (1,0... Number	Scoring Time (1,00... Number
Naive Bayes	0.216	0.020	96	19192	117.845	271.709
Logistic Regression	0.176	0.028	102	6429	141.414	224.090
Decision Tree	0.173	0.032	98	7056	74.074	207.283
Random Forest	0.169	0.022	104	102091	404.040	7050.420

рис. 14: Результат использования AutoModel и сравнение различных моделей по

ключевым метрикам качества

- В разделе Model для дерева решений (Decision Tree) можно увидеть визуализацию, демонстрирующую важность признаков и логику классификации.

### Decision Tree – Model

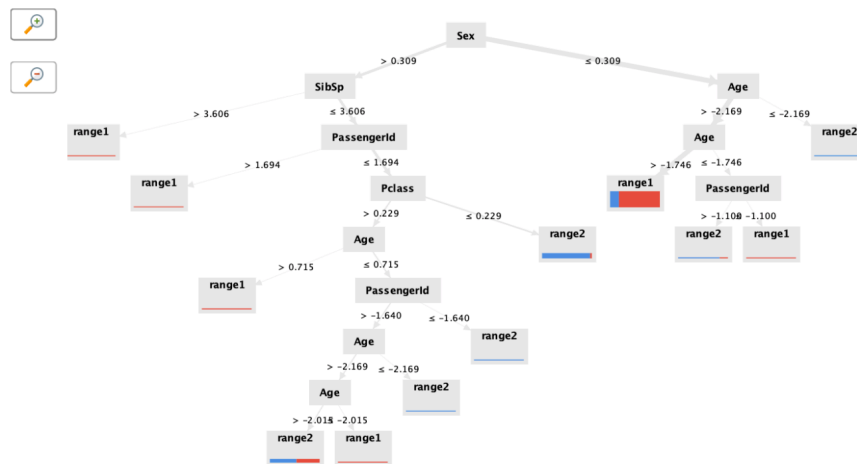


рис. 15: Просмотр архитектуры модели "Дерево решений"

- Также можно рассмотреть другие метрики качества, помимо ассурасы. Например, можно выбрать Recall, F Measure или AUC, которые также будут релевантными показателями.

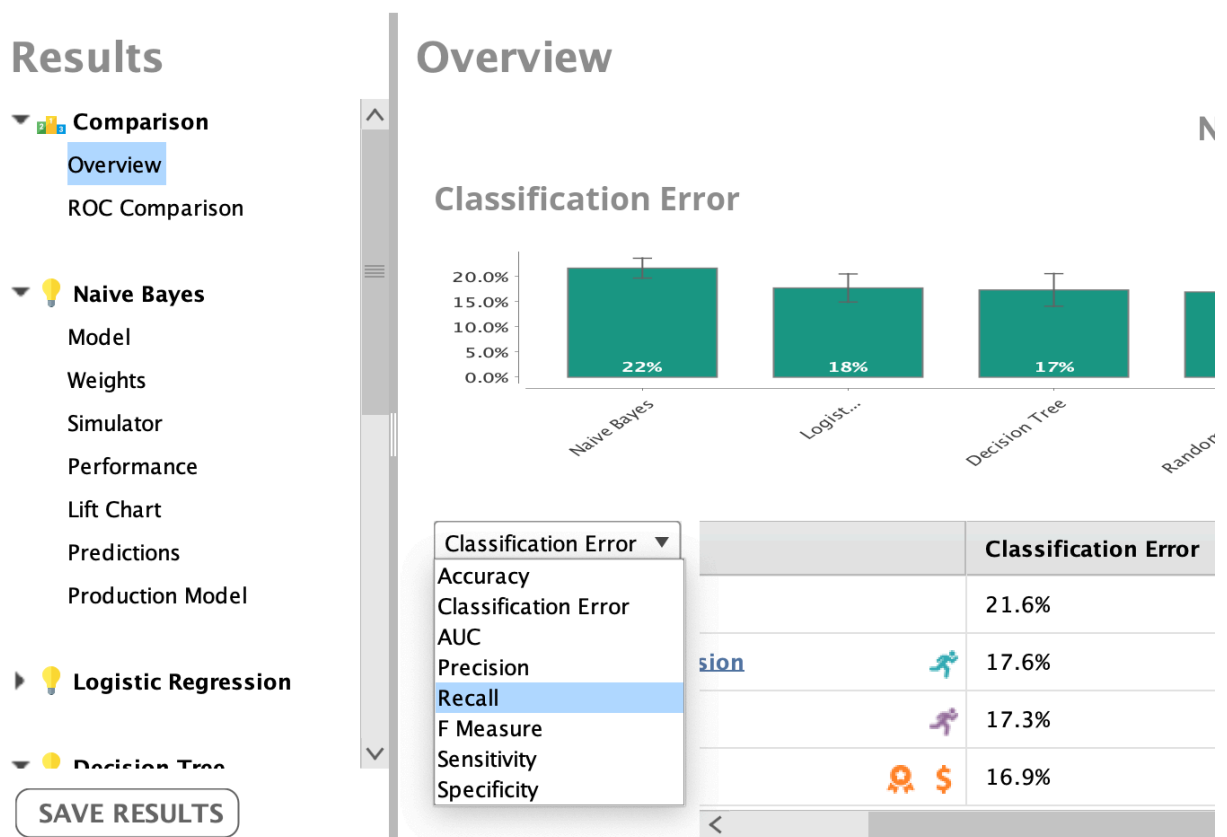


рис. 16: Сравнение моделей по выбранной метрике качества

## Приобретенные навыки

- Умение работать с разными типами графиков в RapidMiner (Altair AI Studio) для исследования структуры данных.
- Опыт настройки параметров визуализации: количество корзин в гистограмме, цветовые схемы, легенды, фильтры.
- Навыки интерпретации результатов моделей с помощью визуальных инструментов (дерево решений, матрица ошибок, метрики точности).
- Способность применять визуальный анализ для принятия обоснованных решений в процессе машинного обучения.

## Обобщенная задача для выполнения индивидуального варианта

Предложенный вам в индивидуальном варианте набор данных будет содержать как числовые, так и категориальные признаки. В ходе выполнения практической работы вам требуется выполнить в RapidMiner (Altair AI Studio) следующие шаги:

Загрузка и предобработка данных:

- Импортируйте данные через оператор Read CSV.
- Проверьте корректность типов столбцов, устраните пропуски и, при необходимости, закодируйте категориальные признаки (Replace / Nominal to Numerical).
- Выполните нормализацию числовых признаков (оператор Normalize).

Изучение структуры данных через визуализацию:

- Постройте гистограмму для одного из ключевых числовых признаков, чтобы оценить распределение.
- Создайте столбчатую диаграмму для одной категориальной переменной и исследуйте ее связь с целевой меткой (если задача кластеризации — без привязки к целевой).
- Постройте круговую диаграмму для распределения классов (или категорий).
- Постройте диаграмму рассеяния по двум числовым признакам, раскрашенную по одной из категорий или по предсказанию простейшей модели.
- Постройте ящик с усами (Box Plot) для числового признака, сгруппированного по категории.
- Постройте тепловую карту корреляций для всех числовых признаков, чтобы оценить их взаимосвязь.

Сравнение моделей:

- Запустите AutoModel для задачи классификации или регрессии, оцените несколько моделей.
- Сравните их точность и другие метрики, визуализируйте структуру лучшей модели (например, дерево решений) и важность признаков.

Интерпретация результатов и выводы:

- На основе визуализаций сформулируйте ключевые наблюдения о распределении данных и влиянии признаков.
- Опишите, какие виды графиков оказались наиболее информативными и почему.

Распределение вариантов





