


Logiном 6.2 - Технический обзор НОВОВВЕДЕНИЙ

Алексей Субботин

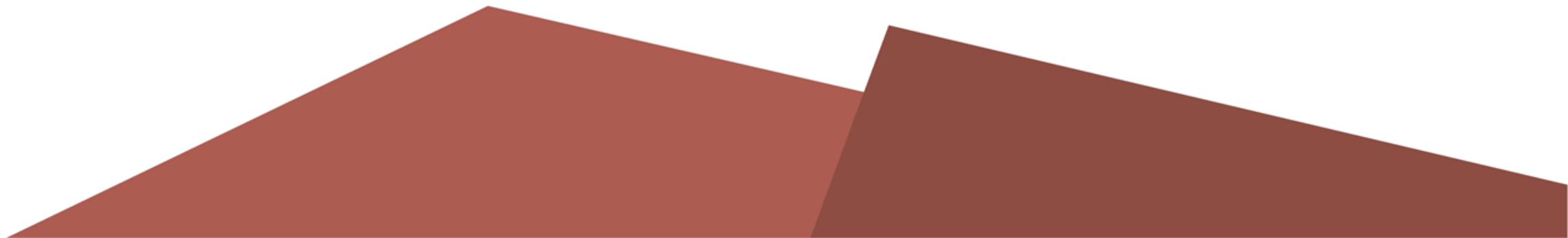
Logiном Company



Новое в версиях 6.1-6.2

1. Востребованные алгоритмы
2. Удобное проектирование
3. Легкость интерпретации

Алгоритмы




Регрессия и бинарная классификация

Логистическая и линейная регрессия –
самые популярные методы:

- Алгоритмы первого выбора
- Широкая применимость
- Понятная интерпретация

Логистическая и линейная регрессия

Автоматическая настройка для «прикладного» моделирования с выбором приоритета:

- Максимальная скорость
 - Повышенная скорость
 - Средние скорость/точность
 - Повышенная точность
 - Максимальная точность
- 

Логистическая и линейная регрессия

Отбор факторов и построение модели с выбором приоритетов достоверные/недостоверные данные, больше/меньше факторов, точность/скорость:

- Классические алгоритмы: **Enter, Forward, Backward, Stepwise**
- Современные алгоритмы с защитой от переобучения: **LASSO, Ridge, Elastic-Net**

Логистическая и
линейная регрессия –
экспертный режим
для тонкой настройки
моделей

Детальная настройка логистической регрессии




Настройки метода	
Точность решения	0,0001
Порог отсеечения	0,5
Включить в модель константу	<input checked="" type="checkbox"/>
Устранение мультиколлинеарности	<input type="checkbox"/>

Настройки расчета статистики	
Рассчитать доверительный интервал	<input type="checkbox"/>
% доверительного интервала	95
Режим расчета статистики	Для финальной модели

Настройки регуляризации	
Установка коэффициента L1-регуляризации	Задать вручную
Коэффициент L1-регуляризации	1
Установка коэффициента L2-регуляризации	Задать вручную
Коэффициент L2-регуляризации	1

Настройки отбора факторов	
Критерий отбора факторов	Информационный критерий Акаике
Порог значимости при добавлении фактора	0,05
Порог значимости при исключении фактора	0,1
Иерархия взаимодействий	Для всех

Логистическая и линейная регрессия

1. Построение модели в пару кликов
 2. Отбор факторов и подбор лучшей модели
 3. Тонкая настройка для экспертов
- 

Сценарий package Визуализаторы

Пакеты Скoring Построение скоринговой карты Сценарий Логистическая регрессия Визуализаторы

Таблица Дерево Информация о модели Нулевые значения Шаги построения Порог значимости 0

Модель	Показатель	Изменение поля	Поля
1.3.20.24.15	34,417238	ab Земельный участок Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Земельный участок Метка
1.3.20.24.16	35,176013	ab Прописка в данном районе Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Прописка в данном районе Метка
1.3.20.24.17	35,200666	ab Гараж Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Гараж Метка
1.3.20.24.18	34,249080	ab Класс предприятия Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Класс предприятия Метка
1.3.20.24.19	27,799777	ab Время работы предприятия Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Время работы предприятия Метка
1.3.20.24.20	26,737878	ab Специализация Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Специализация Метка
1.3.20.24.21	35,178121	ab Должность Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Должность Метка
1.3.20.24.22	26,504717	ab Срок работы на предприятии Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Срок работы на предприятии Метка
1.3.20.24.23	26,476447	ab Срок работы по специальности Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Срок работы по специальности Метка
1.3.20.24.24	10,414540	ab Среднемес. расход Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка
1.3.20.24.24.1	10,416061	ab Сумма кредита Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Сумма кредита Метка
1.3.20.24.24.2	10,416107	ab Стоимость кредита Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Стоимость кредита Метка
1.3.20.24.24.3	6,493915	ab Цель кредитования Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Цель кредитования Метка
1.3.20.24.24.4	2,778093	ab Возраст Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Возраст Метка
1.3.20.24.24.5	2,777982	ab Пол Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Пол Метка
1.3.20.24.24.6	10,347677	ab Образование Метка	ab Срок кредита Метка, ab Отрасль предприятия Метка, ab Среднемес. доход Метка, ab Среднемес. расход Метка, ab Образование Метка

Модель 1.3.20.24.22

Показатель	Значение
Константа	true
-2 Log Likelihood	26,504717
R2	0,820793
R2 корр.	0,793815
Chi-квадрат	-121,395142
Число степеней свободы	93
Значимость	0,000000
AIC	58,504717
AICc	64,482739
BIC	101,418817
HQC	75,904826

ab Давать кредит	Значение	Кол-во
Событие Да	Да	59
Не событие Нет	Нет	90

Атрибут	Коэффициент	Стандартная ошибка	Коэффициент Вальда	Значимость	Отношение шансов	Нижняя граница ДИ	Верхняя граница ДИ
90 Константа	39,114160	0,435014	8 084,666687	0,000000	97,065 298 402 660 35...	40 916 479 447 626 100,0000...	230 2...
ab Срок кредита Метка							
ab 12 ≤... < 18	-35,271012	0,773467	2 079,469099	0,000000	0,000000	0,000000	0,000000
ab 18 ≤... < 24	-38,940739	1,186371	1 121,842	0,000000	0,000000	0,000000	0,000000
ab 24 ≤... < 30	-50,358531	3,128577	259,090982	0,000000	0,000000	0,000000	0,000000
ab от 30	-53,656051	2,455003	477,67...	0,000000	0,000000	0,000000	0,000000
ab Отрасль предприятия Метка							
ab [Издательская деятельность; Иное; Здравоох...	-22,978827	0,654699	1 231,891158	0,000000	0,000000	0,000000	0,000000
ab [Наука и культура; Сельское хозяйство; Рекла...	-5,716854	0,000000	∞	1,000000	0,003290	0,003290	0,003290
ab [Образование (комерч.); Правоохранительны...	-1,584128	0,732123	4,6817...	0,03290	0,047932	0,047932	0,047932
ab [Торговля оптовая, поср. деятельность; Маши...	-3,788231	2,593737	2,133149	0,144144	0,022636	0,000131	0,000131
ab Срок работы на предприятии Метка							
ab 4 ≤... < 5	1,024630	2,375461	0,1860...	0,661111	0,024908	0,024908	0,024908
ab до 3	4,244695	2,134992	3,952761	0,046795	69,734461	1,005037	1,005037
ab от 5	5,828933	0,859189	46,025679	0,000000	339,995541	61,728923	61,728923
ab Среднемес. доход Метка							
ab 4500 ≤... < 5000	-7,966416	1,872467	18,100...	0,000000	0,000000	0,000000	0,000000
ab 5500 ≤... < 6000	-6,070843	0,998461	36,968802	0,000000	0,002309	0,000318	0,000318

Отчет по логистической и линейной регрессии:

- Шаги построения
- Показатели качества
- Детальное описание



Выбор диаграммы

- ROC-кривая
 - PR-кривая
 - Базовые показатели
 - Диаграмма точности
 - Диаграмма равновесия
 - % распознанных событий
 - Диаграмма роста
 - Диаграмма отклика
 - Gain-диаграмма
- Кумулятивная
- 10 диапазонов

- ### Множества
- Обучающее
 - Тестовое

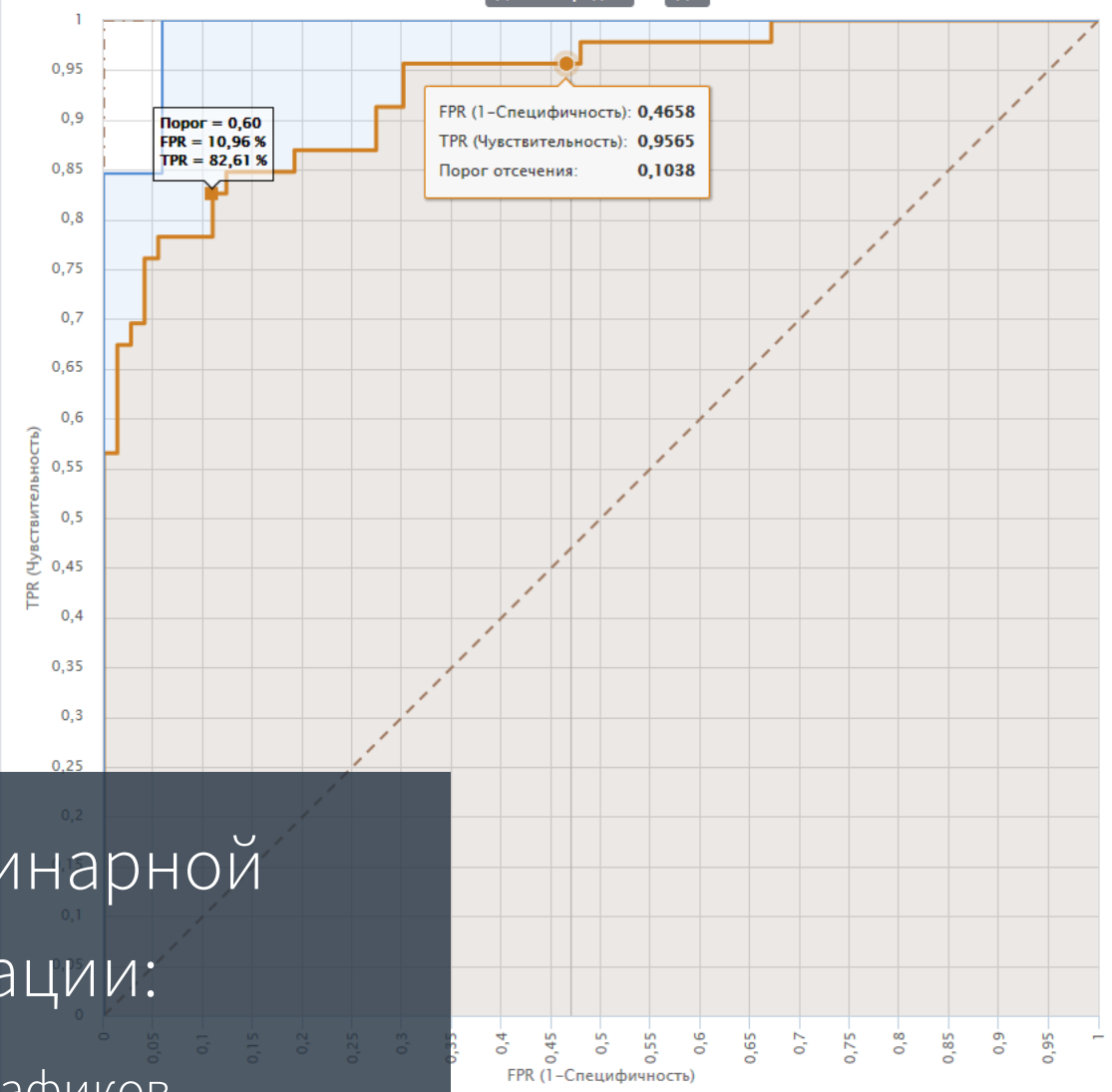
Порог отсечения

Заданный

Значение порога: 0,6

ROC-кривая

Событие: Давать кредит = Да



- Обучающее множество
- Порог отсечения
- Базовая линия
- Идеальная линия
- Тестовое множество
- Порог отсечения

Оценки классификации

Показатель	Множества	
	Обучающее	Тестовое
Оценки классификатора		
AUC ROC	0,9324	0,9910
AUC PR	0,9201	0,9881
Коэффициент Джини	0,8648	0,9819
KS	72,7814	94,1176
Порог отсечения		
Значение	0,6000	0,6000
TPR (Чувствительность)	0,8261	1,0000
TNR (Специфичность)	0,8904	0,9412
FPR (1-Специфичность)	0,1096	0,0588
PPV	0,8261	0,9286
F1 Score	0,8261	0,9630
MCC	0,7165	0,9349

Матрицы ошибок

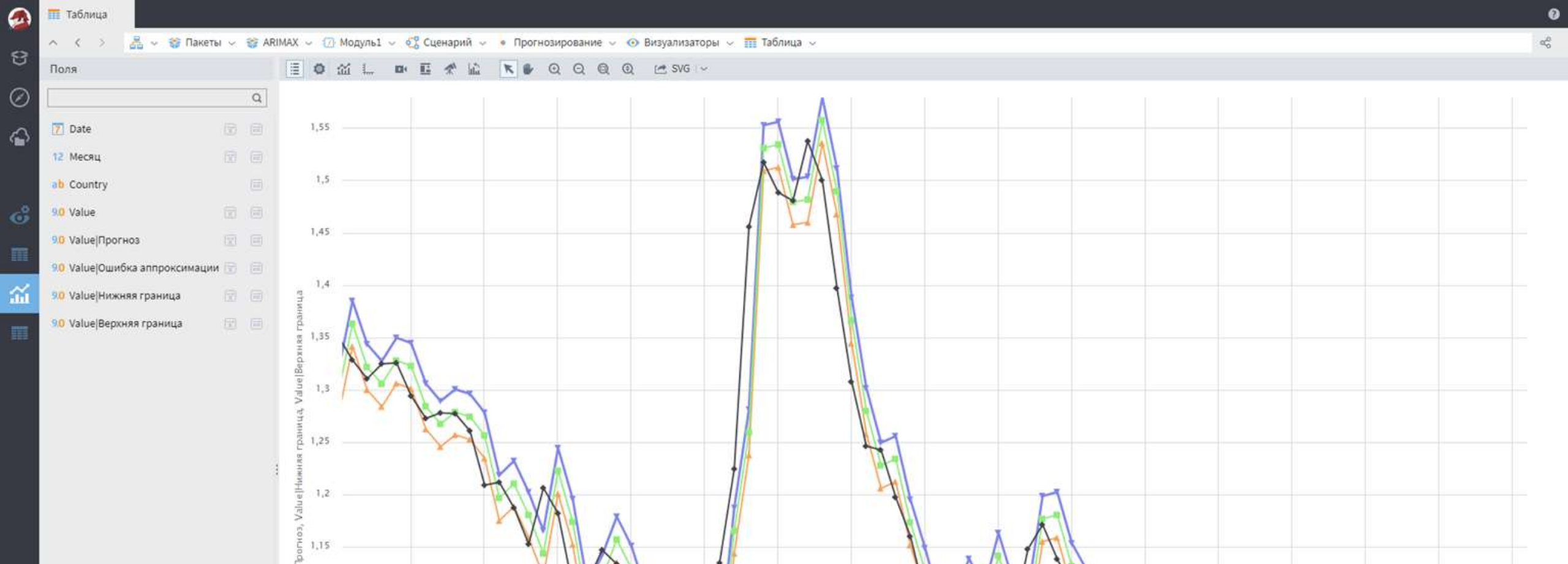
Классифицировано	Фактически		Итого
	Событие	Не-событие	
Обучающее множество	46	73	
Событие	38	8	46
Не-событие	8	65	73
Тестовое множество	13	17	
Событие	13	1	14
Не-событие	0	16	16

Распознано

Обучающее множество	86,55%
Тестовое множество	96,67%

Качество бинарной классификации:

1. Все виды графиков
2. Все известные индикаторы



Прогнозирование ARIMAX

1. Встроенный учет сезонности
2. Включение внешних факторов
3. Автоподбор структуры

→ Value ■ Value|Прогноз ▲ Value|Нижняя граница ◆ Value|Верхняя граница

Настройка конечных классов

Состояние входа: Активировано

Столбец	Индекс F
12 Срок кредита	2,38
90 Сумма кредита	2,26
90 Стоимость кредита	2,26
ab Отрасль предпри...	0,88
12 Время работы пр...	0,74
90 Площадь квартиры	0,41
12 Срок эксплуатац...	0,39
12 Возраст	0,29
ab Основное напра...	0,29
12 Среднемес. расход	0,23

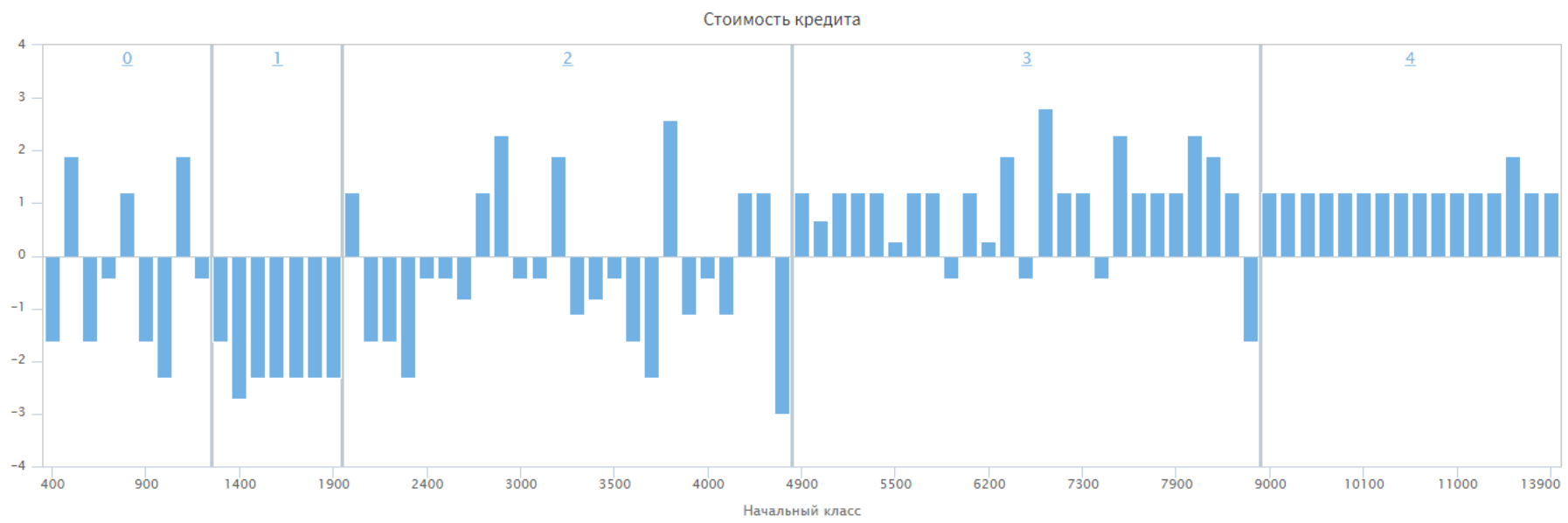
Конечные классы

Минимальный вес, %:

Максимальное кол-во:

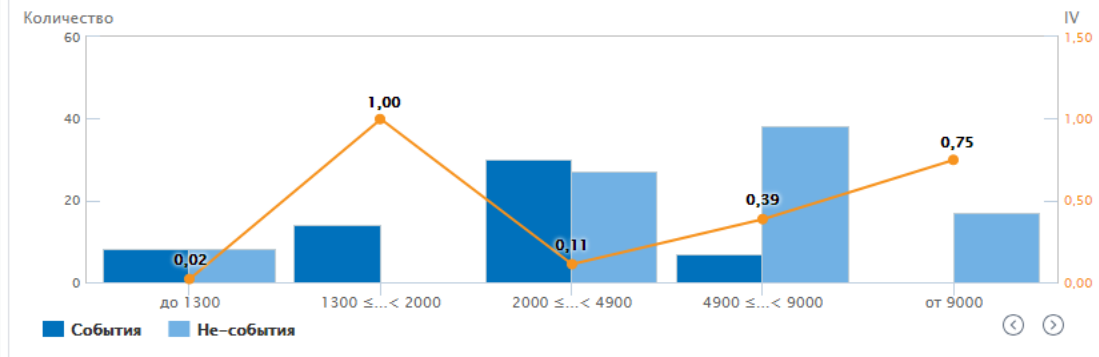
Оптимизация

Равномерность, %:



Конечные классы:

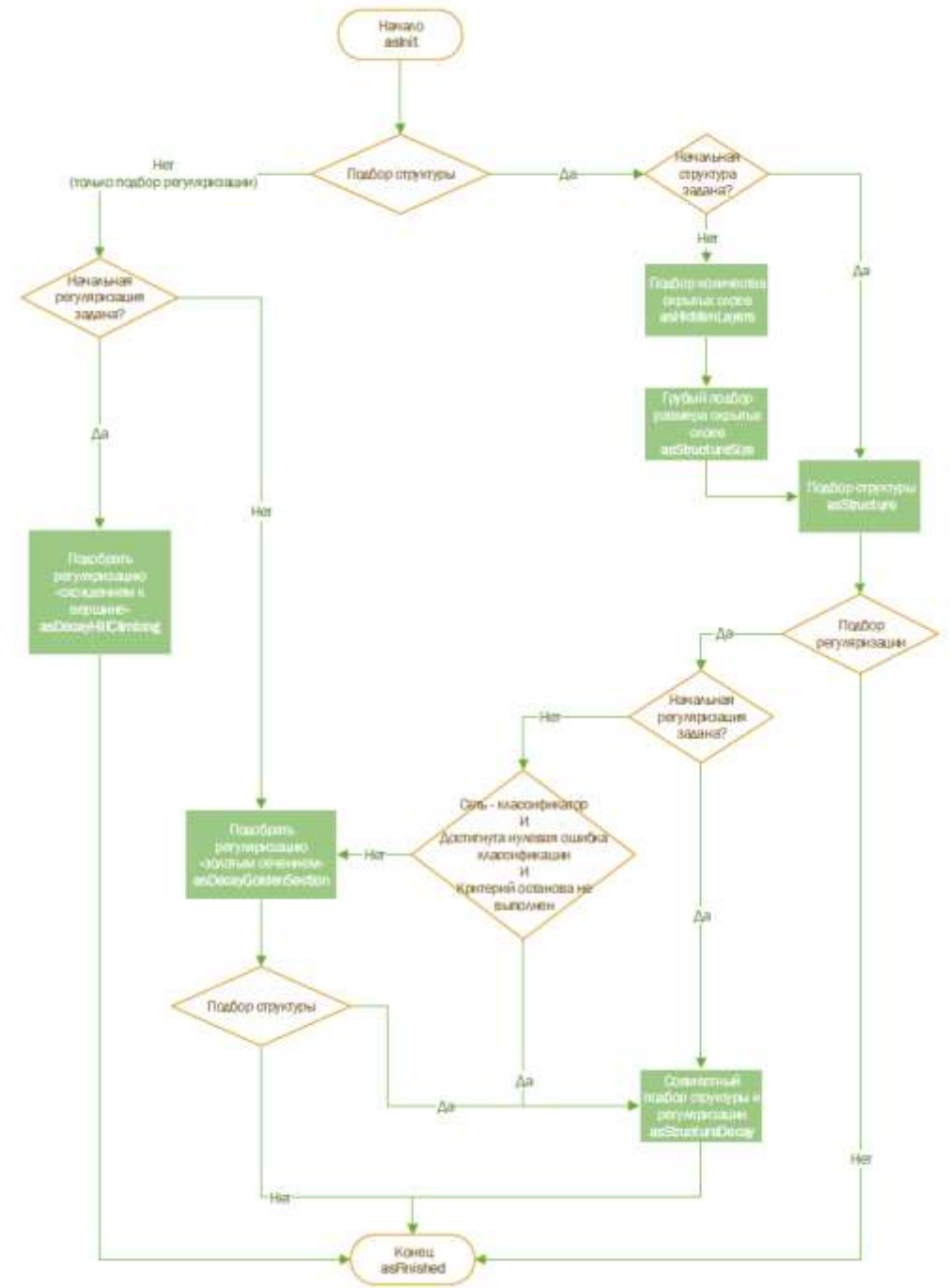
1. Все показатели на одном экране
2. Ручная правка разбиения
3. Поддержка внешних диапазонов



Нейросети

Автоподбор структуры и степени регуляризации:

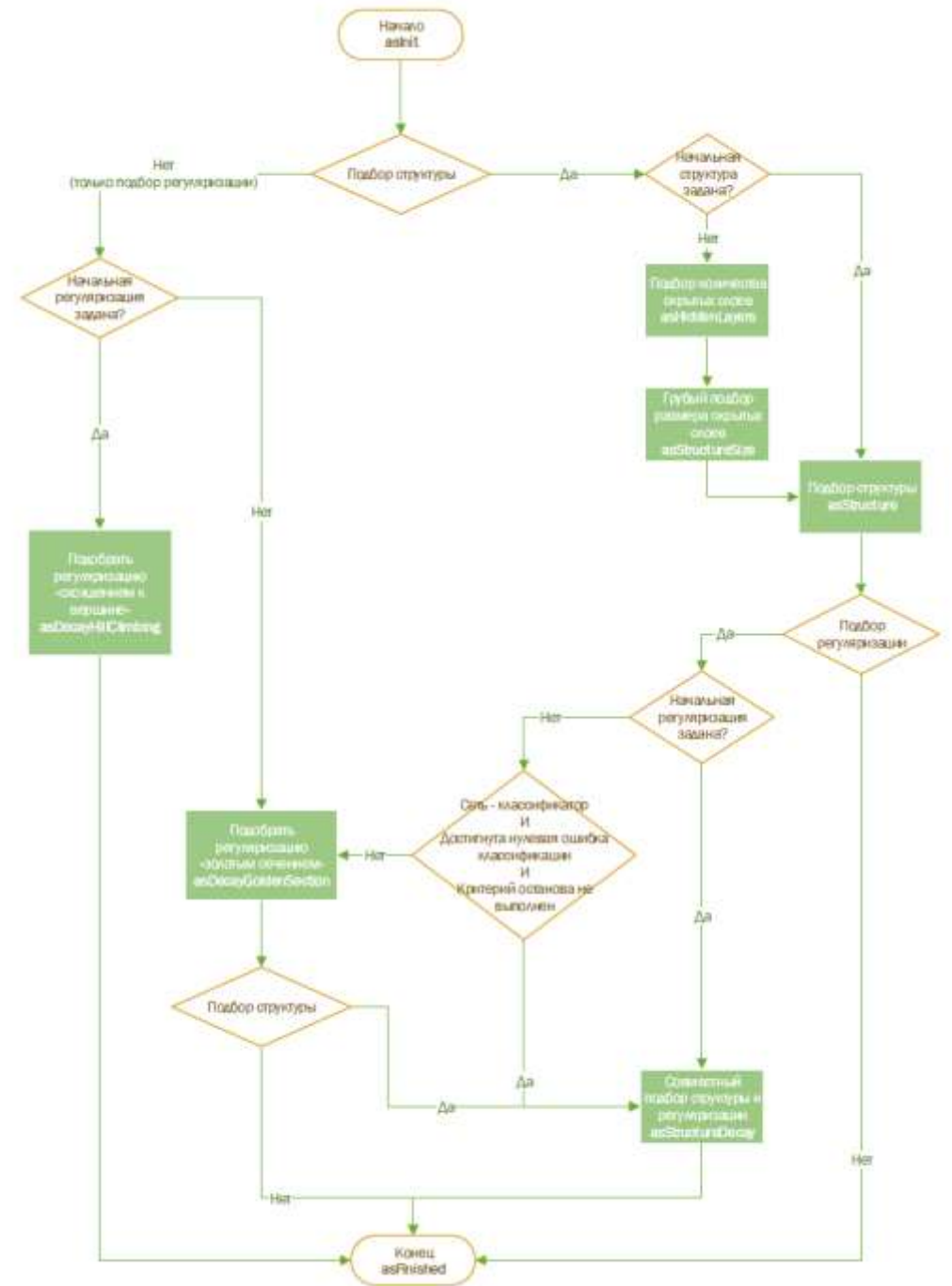
- Оптимальная структура
- Защита от переобучения
- Экономия времени и памяти




Нейросети

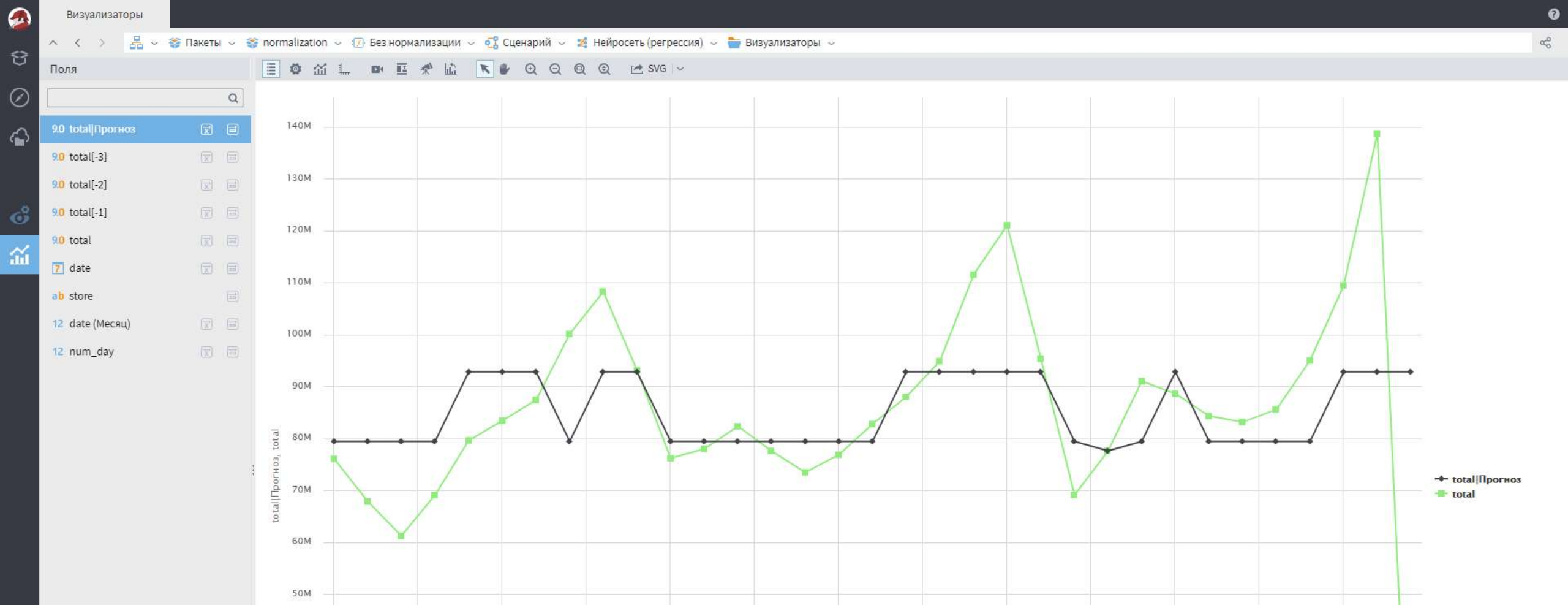
Метаалгоритм, не требующий участия пользователя с начальной точкой и без:

- Golden Section
- Hill Climbing
- Exhaustive Search



Отсутствие нормализации
(денормализации) усложняет
работу аналитика:

- Плохо обусловленные данные
 - Долгое сходжение алгоритмов
 - Неадекватные результаты
- 



Модель прогнозирования
без нормализации данных

Сложности ручной нормализации

- Самостоятельный расчет, сохранение, поддержание и синхронизация настроек
- Поддержание списков уникальных значений
- Трудности при замене входных данных
- Цена ошибки – неработающая модель

Нормализация/денормализация

Непрерывные данные

1. Нормализация MIN-MAX
2. Нормализация [0; 1]
3. Нормализация [-1; 1]
4. Абсолютное масштабирование
5. Стандартизация
6. Отношение

Дискретные данные

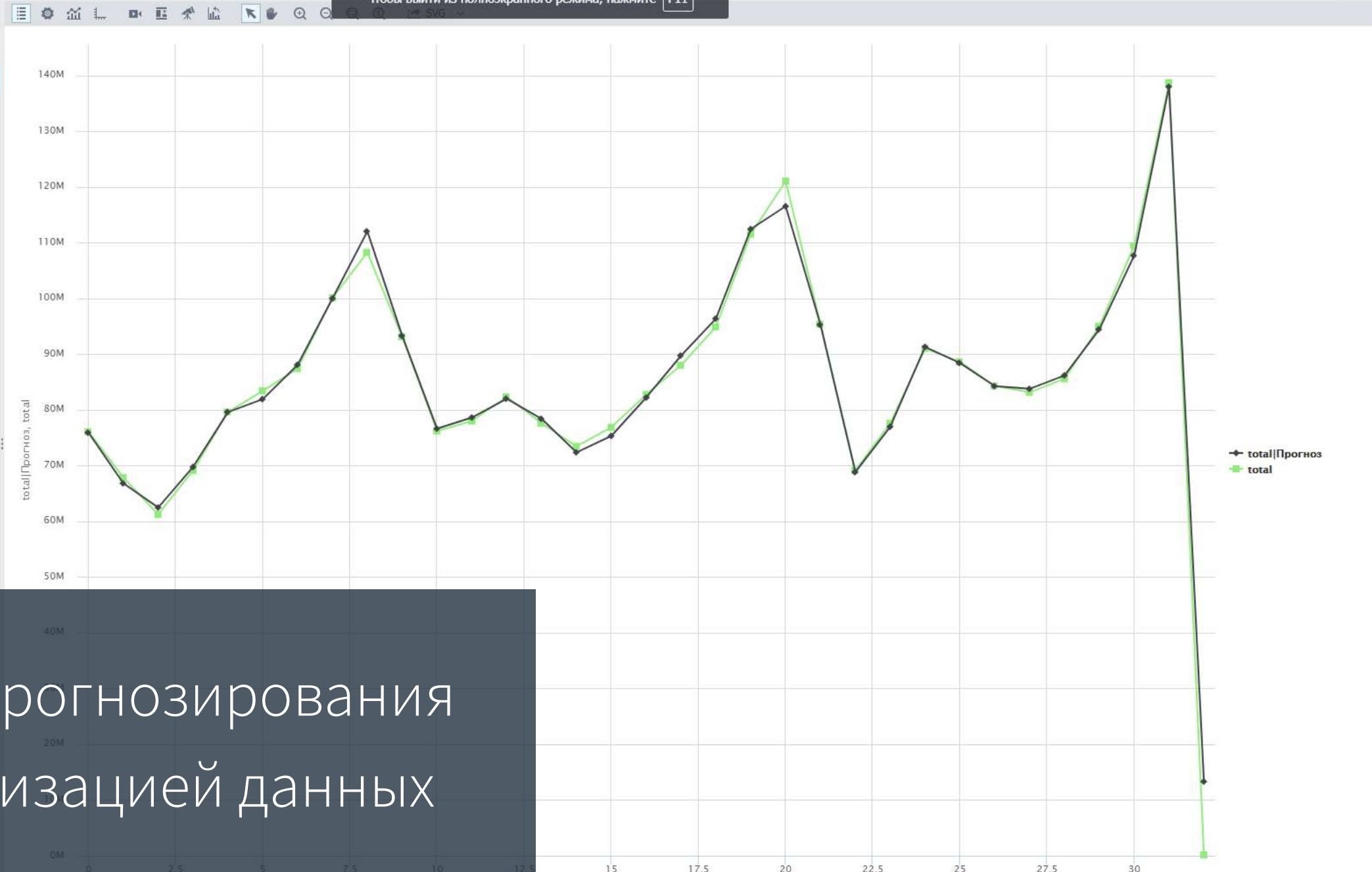
1. Индикатор
2. Индикатор без опорной категории
3. Отклонение
4. Простой
5. Разность
6. Обратная разность
7. Helmert
8. Обратный Helmert
9. Индекс уникального значения

Нормализация/денормализация встроена в Data Mining алгоритмы:

- Логистическая регрессия
- Линейная регрессия
- Нейросеть (классификация, регрессия)
- Кластеризация (k-means)
- EM-кластеризация
- SONN
- Прогнозирование (ARIMAX)

Поля

- 9.0 total|Прогноз
- 9.0 total нормализованное|Прогноз
- 9.0 total[-3] нормализованное
- 9.0 total[-2] нормализованное
- 9.0 total[-1] нормализованное
- 9.0 total нормализованное
- 9.0 total[-3]
- 9.0 total[-2]
- 9.0 total[-1]
- 9.0 total
- 7 date
- ab store
- 12 date (Месяц)
- 12 num_day



Модель прогнозирования с нормализацией данных

Обобщающая способность моделей

Разбиение на множества встроено в компоненты там, где оно необходимо:

- Случайный и последовательный сэмплинг
- Автоматический расчет всех показателей на обучающем и тестовом множествах

Обобщающая способность моделей

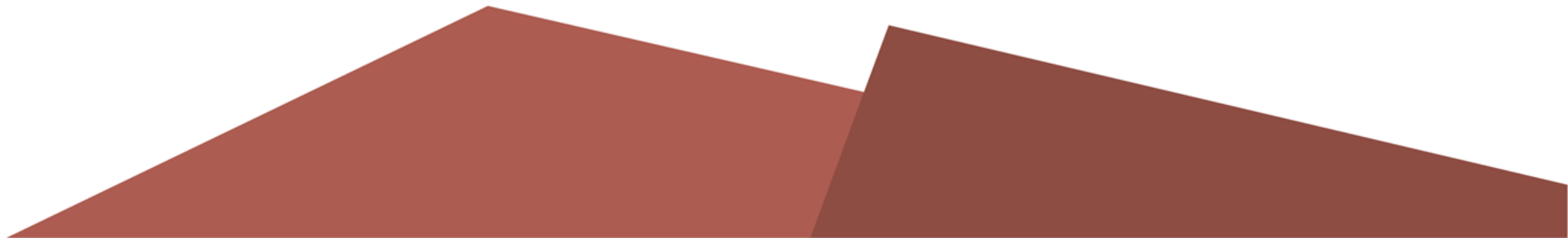
Кросс-валидация встроена в компоненты, допускающие подобную проверку:

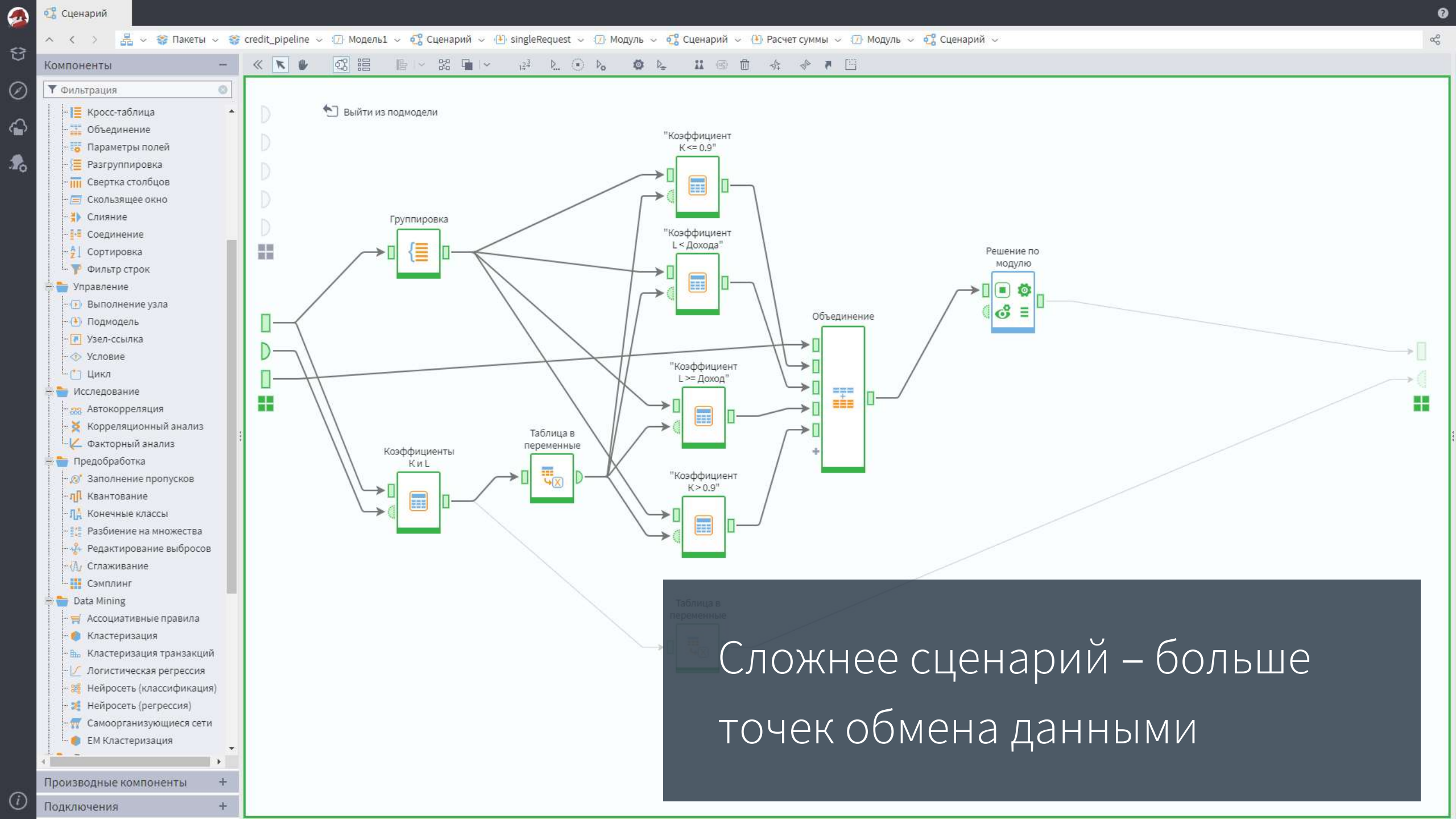
- K-folds
- Shuffle split

Обобщающая способность моделей

- Тестирование и кросс-валидация без дополнительных усилий
- Выходные показатели унифицированы среди различных компонентов
- Удобное сравнение моделей различной природы

Удобство проектирования



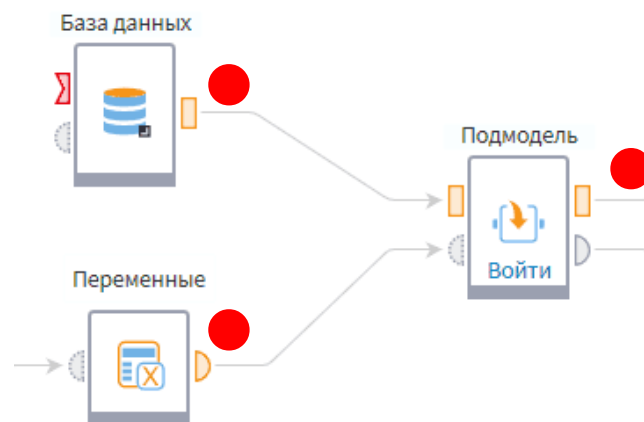
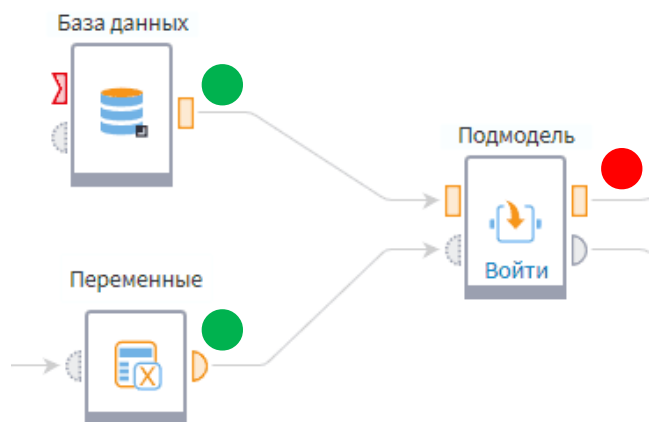
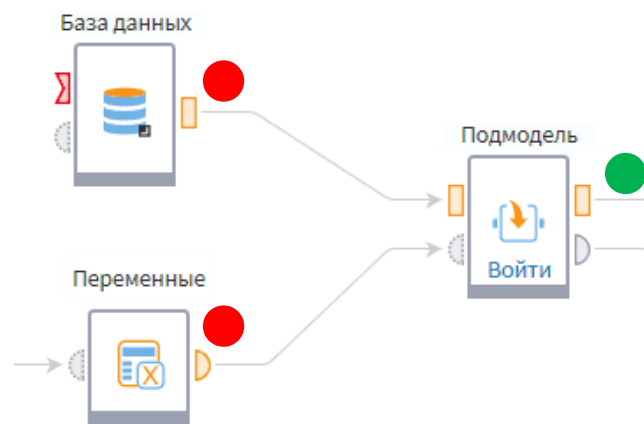
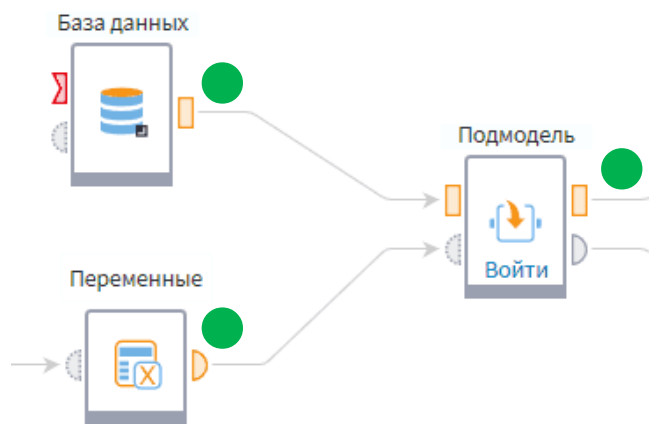


- Кросс-таблица
- Объединение
- Параметры полей
- Разгруппировка
- Свертка столбцов
- Скользящее окно
- Слияние
- Соединение
- Сортировка
- Фильтр строк
- Управление
- Выполнение узла
- Подмодель
- Узел-ссылка
- Условие
- Цикл
- Исследование
- Автокорреляция
- Корреляционный анализ
- Факторный анализ
- Предобработка
- Заполнение пропусков
- Квантование
- Конечные классы
- Разбиение на множества
- Редактирование выбросов
- Сглаживание
- Сэмплинг
- Data Mining
- Ассоциативные правила
- Кластеризация
- Кластеризация транзакций
- Логистическая регрессия
- Нейросеть (классификация)
- Нейросеть (регрессия)
- Самоорганизующиеся сети
- EM Кластеризация

Таблица в переменные

Сложнее сценарий – больше точек обмена данными

Автосинхронизация: кейсы




Автосинхронизация


Обязательные поля:

- Настроенные пользователем
- С заданным назначением

Необязательные поля:

- Пробрасываем насквозь
 - Данные – только по требованию
- 

Автосинхронизация

- Минимум действий пользователя при подмене источников данных
 - Все конфликты подсвечиваем визуально
 - Покажем мастер, только если потребуется вмешательство пользователя
- 

Эволюция Logiном

1. Больше алгоритмов анализа
2. Повышение качества анализа
3. Автоматизация рутинной работы
4. Улучшение интерпретируемости

loginom.ru