

**Второе издание,  
обновленное  
и дополненное**



**DMK**  
ИЗДАТЕЛЬСТВО

ИИ «Гевисста»  
**G**

**Артем Груздев**

# Прогнозное моделирование в IBM SPSS Statistics, R и Python

**Метод деревьев решений  
и случайный лес**

**УДК 519.7:004.9IBM SPSS Statistics**  
**ББК 21.18с**  
**Г90**

**Груздев А. В.**

**Г90** Прогнозное моделирование в IBM SPSS Statistics, R и Python: метод деревьев решений и случайный лес. – М.: ДМК Пресс, 2018. – 642 с.: ил.

**ISBN 978-5-97060-539-4**

Данная книга представляет собой практическое руководство по применению метода деревьев решений и случайного леса для задач сегментации, классификации и прогнозирования. Каждый раздел книги сопровождается практическим примером. Кроме того, книга содержит программный код SPSS Syntax, R и Python, позволяющий полностью автоматизировать процесс построения прогнозных моделей. Автором обобщены лучшие практики использования деревьев решений и случайного леса от таких компаний, как Citibank N.A., Transunion и DBS Bank.

Издание будет интересно маркетологам, риск-аналитикам и другим специалистам, занимающимся разработкой и внедрением прогнозных моделей.

**УДК 519.7:004.9IBM SPSS Statistics**  
**ББК 21.18с**

Все права защищены. Любая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами без письменного разрешения владельцев авторских прав.

Материал, изложенный в данной книге, многократно проверен. Но поскольку вероятность технических ошибок все равно существует, издательство не может гарантировать абсолютную точность и правильность приводимых сведений. В связи с этим издательство не несет ответственности за возможные ошибки, связанные с использованием книги.

ISBN 978-5-97060-539-4

© Груздев А. В., 2018  
© Оформление, издание, ДМК Пресс, 2018

# Содержание

<b>От рецензента</b> .....	10
<b>Предисловие</b> .....	11
<b>Глава 1. Введение в метод деревьев решений</b> .....	14
1.1. Введение в методологию деревьев решений .....	14
1.2. Преимущества и недостатки деревьев решений .....	19
1.3. Задачи, выполняемые с помощью деревьев решений .....	20
Вопросы к главе 1 .....	22
<b>Часть I. ПОСТРОЕНИЕ ДЕРЕВЬЕВ РЕШЕНИЙ И СЛУЧАЙНОГО ЛЕСА В IBM SPSS STATISTICS</b> .....	23
<b>Глава 2. Основы прогнозного моделирования с помощью деревьев решений CHAID</b> .....	24
2.1. Запуск процедуры Деревья классификации .....	24
2.2. Четыре метода деревьев решений .....	26
2.3. Шкалы переменных .....	29
2.4. Определение необходимого размера выборки .....	31
2.5. Знакомство с методом CHAID .....	32
2.5.1. Описание алгоритма .....	32
2.5.2. Немного о тесте хи-квадрат .....	35
2.5.3. Немного об F-тесте .....	37
2.5.4. Способы объединения категорий предикторов .....	38
2.5.5. Поправка Бонферрони .....	38
2.5.6. Иллюстрация работы CHAID на конкретном примере .....	38
2.6. Построение и интерпретация дерева классификации CHAID .....	43
2.6.1. Сводка для модели .....	45
2.6.2. Диаграмма дерева .....	46
2.6.3. Выигрыши для узлов .....	48
2.6.4. Таблицы классификации и риска .....	49
2.7. Работа с прогнозами модели .....	51
2.7.1. Получение результатов классификации .....	51
2.7.2. Сохранение прогнозов модели в файле данных .....	52
2.7.3. Самостоятельное построение таблицы классификации и изменение порогового значения вероятности .....	57
2.8. Анализ ROC-кривой .....	66
2.8.1. Терминология анализа ROC-кривой .....	66
2.8.2. Оценка дискриминирующей способности модели и выбор порогового значения с помощью ROC-кривой .....	73
2.9. Диагностика качества модели .....	79
2.9.1. Обобщающая способность, переобучение и недообучение .....	79
2.9.2. Методы проверки модели .....	80
2.9.3. Общие правила интерпретации результатов проверки .....	83
2.9.4. Методы проверки модели, реализованные в процедуре Деревья классификации .....	85

2.9.5. Практическое применение методов проверки в процедуре Деревья классификации .....	86
2.9.6. Самостоятельное разбиение набора данных на обучающую и контрольную выборки для осуществления проверки .....	97
2.10. Дополнительные настройки вывода результатов .....	101
2.10.1. Настройки вывода дерева .....	101
2.10.2. Построение таблицы дерева .....	102
2.10.3. Настройки вывода статистик .....	103
2.10.4. Построение таблиц выигрышей для узлов и процентилей.....	105
2.10.5. Настройки вывода графиков .....	107
2.10.6. Построение графиков выигрышей, индексов и откликов.....	109
2.10.7. Настройки вывода правил классификации.....	111
2.10.8. Применение правил классификации к новому набору данных .....	112
2.11. Построение дерева регрессии CHAID.....	122
2.12. Использование принудительной переменной расщепления.....	127
Выводы и рекомендации .....	129
Вопросы к главе 2.....	130

### **Глава 3. Продвинутое моделирование**

<b>с помощью деревьев решений CHAID.....</b>	<b>133</b>
3.1. Построение деревьев CHAID с измененными критериями.....	133
3.1.1. Настройка правил остановки .....	133
3.1.2. Построение деревьев CHAID с измененными правилами остановки .....	134
3.1.3. Настройка статистических тестов для разбиения узлов и объединения категорий предикторов.....	140
3.1.4. Построение дерева CHAID с измененными статистическими тестами.....	141
3.1.5. Настройка обработки количественных предикторов .....	142
3.1.6. Построение дерева CHAID с измененным числом интервалов для количественных предикторов .....	143
3.2. Метод Исчерпывающий CHAID .....	144
3.3. Обзор параметров деревьев решений.....	145
3.4. Работа с пропусками в методе CHAID .....	147
3.4.1. Настройка обработки пропущенных значений.....	147
3.4.2. Построение дерева CHAID на основе данных, содержащих пропуски .....	150
3.5. Работа со стоимостями ошибочной классификации в методе CHAID.....	151
3.5.1. Настройка стоимостей ошибочной классификации .....	151
3.5.2. Построение дерева CHAID с измененными стоимостями ошибочной классификации .....	154
3.6. Работа с прибылями в методе CHAID.....	157
3.6.1. Настройка прибылей .....	157
3.6.2. Построение дерева CHAID с заданными значениями прибыли.....	158
3.7. Работа со значениями .....	162
3.8. Применение метода CHAID для биннинга переменных (на примере конкурсной задачи ОТП Банка).....	165
3.8.1. Преимущества и недостатки биннинга .....	165
3.8.2. Предварительная подготовка данных .....	167
3.8.3. Определение важности переменных с помощью случайного леса.....	184
3.8.4. Анализ мультиколлинеарности .....	187
3.8.5. Выполнение биннинга переменных на основе CHAID.....	188

3.8.6. Построение моделей логистической регрессии на основе исходных предикторов и предикторов, категоризированных с помощью CHAID .....	194
3.8.7. Выполнение биннинга переменных с помощью процедуры Оптимальная категоризация .....	199
3.8.8. Построение модели логистической регрессии на основе оптимально категоризированных предикторов.....	202
3.8.9. Преобразование количественных переменных для максимизации нормальности .....	203
3.8.10. Построение модели логистической регрессии с использованием CHAID и преобразования корня третьей степени.....	207
3.9. Построение ансамбля логистической регрессии и дерева CHAID (на примере конкурсной задачи Tinkoff Data Science Challenge).....	208
Выводы и рекомендации .....	218
Вопросы к главе 3.....	219
<b>Глава 4. Построение деревьев решений CRT и QUEST .....</b>	<b>220</b>
4.1. Знакомство с методом CRT .....	220
4.1.1. Описание алгоритма .....	221
4.1.2. Мера Джини .....	222
4.1.3. Внутриузловая дисперсия .....	223
4.1.4. Метод отсеечения ветвей на основе меры стоимости-сложности .....	224
4.1.5. Обработка пропущенных значений.....	225
4.1.6. Иллюстрация работы CRT на конкретном примере.....	225
4.2. Построение дерева классификации CRT.....	228
4.3. Построение дерева CRT с измененными критериями .....	231
4.3.1. Настройка мер неоднородности для отбора предикторов и расщепления узлов .....	232
4.3.2. Настройка отсеечения ветвей.....	233
4.3.3. Построение дерева CRT с последующим отсечением ветвей .....	234
4.3.4. Настройка суррогатов для обработки пропущенных значений .....	235
4.3.5. Построение дерева CRT на основе данных, содержащих пропуски .....	236
4.4. Вывод важности предикторов.....	239
4.5. Работа с априорными вероятностями в методе CRT.....	240
4.5.1. Настройка априорных вероятностей .....	240
4.5.2. Построение дерева CRT с измененными априорными вероятностями .....	241
4.6. Знакомство с методом QUEST.....	243
4.6.1. Описание алгоритма .....	244
4.6.2. Метод отсеечения ветвей на основе меры стоимости-сложности .....	246
4.7. Построение дерева классификации QUEST .....	246
4.8. Сравнение метода QUEST с другими методами деревьев решений .....	248
4.9. Построение дерева QUEST с измененными критериями.....	249
4.9.1. Настройка статистических тестов для отбора предикторов.....	250
4.9.2. Построение дерева QUEST с последующим отсечением ветвей.....	250
Выводы и рекомендации .....	252
Вопросы к главе 4.....	252
<b>Глава 5. Редактор дерева .....</b>	<b>254</b>
5.1. Просмотр диаграммы дерева в Редакторе .....	254
5.2. Просмотр содержимого узла в Редакторе.....	255

5.3. Настройка внешнего вида диаграммы дерева в Редакторе.....	256
5.4. Изменение ориентации диаграммы дерева в Редакторе.....	257
5.5. Настройка содержимого узла в Редакторе.....	257
5.6. Отбор наблюдений в Редакторе.....	258
5.7. Иллюстрация работы в Редакторе дерева на конкретном примере .....	259

## **Глава 6. Построение случайного леса .....**

263

6.1. Введение в методологию случайного леса.....	263
6.1.1. Описание метода.....	263
6.1.2. Оценка качества модели.....	267
6.1.3. Настройка параметров случайного леса .....	270
6.1.4. Важность предикторов.....	271
6.1.5. Графики частной зависимости .....	273
6.1.6. Матрица близостей .....	275
6.1.7. Обработка пропущенных значений.....	276
6.1.8. Обнаружение выбросов.....	276
6.1.9. Преимущества и недостатки случайного леса.....	277
6.1.10. История создания метода.....	278
6.2. Знакомство с процедурой Оценка RanFor .....	278
6.3. Построение ансамбля деревьев классификации .....	282
6.4. Интерпретация результатов, полученных с помощью ансамбля деревьев классификации .....	286
6.4.1. Сводка для модели .....	286
6.4.2. Важность переменных.....	288
6.4.3 Частота использования переменных .....	288
6.4.4 Матрица ошибок прогнозов.....	289
6.4.5. График частоты ошибок.....	290
6.4.6. График важности переменных.....	291
6.4.7. Графики частной зависимости .....	291
6.4.8 Работа с набором прогнозов.....	294
6.5. Проверка построенного ансамбля деревьев классификации на контрольной выборке и применение его к новым данным с помощью процедуры Прогноз RanFor.....	297
6.6. Построение ансамбля деревьев регрессии и интерпретация полученных результатов.....	303
6.7. Проверка построенного ансамбля деревьев регрессии на контрольной выборке и применение его к новым данным с помощью процедуры Прогноз RanFor .....	311
Выводы и рекомендации .....	315
Вопросы к главе 6.....	315

## **Часть II. ПОСТРОЕНИЕ ДЕРЕВЬЕВ РЕШЕНИЙ И СЛУЧАЙНОГО ЛЕСА В R И RYTHON.....**

318

### **Глава 7. Построение деревьев решений CHAID с помощью пакета R CHAID.....**

319

7.1. Построение и интерпретация дерева классификации CHAID .....	319
7.1.1. Подготовка данных .....	319
7.1.2. Построение модели и работа с диаграммой дерева.....	321

7.1.3. Вычисление вероятностей классов и выбор оптимального порога.....	323
7.1.4. Получение спрогнозированных классов зависимой переменной.....	328
7.1.5. Сохранение прогнозов .....	329
7.1.6. Применение модели к новым данным .....	329
7.1.7. Проверка модели.....	330
7.2. Биннинг переменных .....	335
7.2.1. Биннинг в пакете rattle .....	335
7.2.2. Биннинг в пакете smbinning.....	337
Выводы и рекомендации .....	344
Вопросы к главе 7.....	345

## **Глава 8. Построение деревьев решений CRT**

<b>с помощью пакета R rpart</b> .....	346
8.1. Метод отсеечения ветвей на основе стоимости-сложности с кросс-проверкой.....	346
8.2. Построение и интерпретация дерева классификации CRT .....	347
8.2.1. Подготовка данных .....	347
8.2.2. Построение модели и работа с диаграммой дерева.....	348
8.2.3. Прунинг дерева CRT .....	354
8.2.4. Вычисление вероятностей классов.....	356
8.2.5. Построение ROC-кривой и вычисление более точных оценок дискриминирующей способности.....	356
8.2.6. Сохранение спрогнозированных вероятностей .....	359
8.2.7. Применение модели к новым данным .....	359
8.3. Построение и интерпретация дерева регрессии CRT.....	361
8.3.1. Подготовка данных .....	361
8.3.2. Построение модели и работа с диаграммой дерева.....	362
Выводы и рекомендации .....	365
Вопросы к главе 8.....	365

## **Глава 9. Построение случайного леса с помощью пакета R randomForest**

.....	367
9.1. Построение ансамбля деревьев классификации .....	367
9.1.1. Подготовка данных .....	367
9.1.2. Построение модели и получение ООВ-оценки качества.....	369
9.1.3. Важности предикторов .....	374
9.1.4. Графики частной зависимости .....	375
9.1.5. Вычисление вероятностей классов.....	379
9.1.6. Оценка дискриминирующей способности модели с помощью ROC-кривой .....	380
9.1.7. Получение спрогнозированных классов зависимой переменной.....	383
9.1.8. График зазора прогнозов .....	385
9.2. Построение ансамбля деревьев регрессии.....	386
9.2.1. Подготовка данных .....	386
9.2.2. Построение модели и получение ООВ оценки качества.....	387
9.2.3. Важности предикторов .....	388
9.2.4. Графики частной зависимости .....	389
9.2.5. Работа с прогнозами и вычисление среднеквадратичной ошибки.....	391
9.2.6. Улучшение качества прогнозов.....	392
9.2.7. Вычисление коэффициента детерминации .....	393

9.2.8. Получение более развернутого вывода о качестве модели .....	394
9.3. Поиск оптимальных параметров случайного леса с помощью пакета caret .....	395
9.3.1. Схема оптимизации параметров, реализованная в пакете caret .....	395
9.3.2. Настройка условий оптимизации .....	396
9.3.3. Поиск оптимальных параметров для задачи регрессии .....	398
9.3.4. Поиск оптимальных параметров для задачи классификации .....	400
Выводы и рекомендации .....	410
<b>Глава 10. Построение случайного леса с помощью пакета R ranger .....</b>	<b>411</b>
10.1. Построение ансамбля деревьев классификации .....	411
10.2. Построение случайного леса вероятностей .....	433
10.3. Построение случайного леса выживаемости .....	442
Выводы и рекомендации .....	449
<b>Глава 11. Построение распределенного случайного леса с помощью пакета R h2o .....</b>	<b>450</b>
11.1. Решение задачи классификации .....	450
11.1.1. Подготовка данных .....	450
11.1.2. Построение модели и работа с результатами .....	455
11.1.3. Сохранение модели и применение к новым данным .....	466
11.1.4. Поиск оптимальных значений параметров с помощью решетчатого поиска .....	467
11.2. Решение задачи регрессии .....	478
Выводы и рекомендации .....	482
<b>Глава 12. Построение случайного леса в Python .....</b>	<b>483</b>
12.1. Знакомство с Python .....	483
12.1.1. Обзор основных инструментов Python, предназначенных для подготовки и анализа данных .....	483
12.1.2. Беспроблемная работа с программным кодом .....	490
12.2. Построение модели случайного леса и работа с полученными результатами .....	490
12.2.1. Подготовка данных в pandas .....	491
12.2.2. Параметры случайного леса и подгонка модели .....	500
12.2.3. Важности предикторов .....	505
12.2.4. Прогнозы модели и матрица ошибок .....	508
12.2.5. Отчет о результатах классификации: точность, полнота и F-мера .....	509
12.2.6. Построение ROC-кривой и выбор оптимального порога .....	511
12.2.7. Сравнение модели случайного леса с моделью дерева решений .....	514
12.3. Улучшение качества модели случайного леса .....	520
12.3.1. Методы перекрестной проверки, реализованные в scikit-learn .....	520
12.3.2. Поиск оптимальных параметров случайного леса .....	522
12.4. Построение распределенного случайного леса с помощью модуля H2O .....	541
12.4.1. Подготовка данных для построения стандартной модели случайного леса .....	541
12.4.2. Построение стандартной модели случайного леса .....	552
12.4.3. Применение стандартной модели случайного леса к новым данным .....	557
12.4.4. Подготовка данных для моделирования в H2O .....	560
12.4.5. Построение модели случайного леса с помощью класса H2ORandomForestEstimator .....	564



---

12.4.6. Сохранение модели случайного леса, построенной с помощью класса H2ORandomForestEstimator, и применение к новым данным .....	579
12.4.7. Улучшение качества моделей классов RandomForestClassifier и H2ORandomForestEstimator с помощью конструирования новых признаков .....	581
12.4.8. Выполнение решетчатого поиска с помощью класса H2OGridSearch .....	585
12.4.9. Улучшение качества модели с помощью стекинга .....	590
Выводы и рекомендации .....	598
<b>Приложение 1. Предварительная подготовка данных в Python с помощью библиотеки pandas .....</b>	<b>599</b>
<b>Приложение 2. Предварительная подготовка данных в R .....</b>	<b>604</b>
<b>Приложение 3. Визуализация данных в Python с помощью библиотек matplotlib, seaborn и plotly .....</b>	<b>612</b>
<b>Приложение 4. Построение ROC-кривой и вычисление AUC вручную .....</b>	<b>616</b>
<b>Приложение 5. Декомпозиция прогнозов дерева решений и случайного леса с помощью питоновского пакета treeinterpreter для улучшения интерпретабельности .....</b>	<b>622</b>
<b>Ключи к вопросам .....</b>	<b>630</b>
<b>Библиографический список .....</b>	<b>631</b>
<b>Предметный указатель .....</b>	<b>633</b>

## Основы прогнозного моделирования с помощью деревьев решений CHAID

### 2.1. Запуск процедуры *Дерева классификации*

Дерево решений в IBM SPSS Statistics можно построить с помощью процедуры **Дерева классификации**. Для вызова процедуры **Дерева классификации** необходимо в меню **Анализ** выбрать **Классификация** ⇒ **Дерева классификации**.

Вы оказываетесь в главном диалоговом окне **Дерево решений** (рис. 2.1). В поле **Зависимая переменная** необходимо перенести одну зависимую переменную. Кнопка **Категории** позволяет включить/исключить из анализа категории зависимой переменной или задать их как целевые. Указание одной или нескольких категорий как целевых не влияет на модель дерева, оценки рисков и результаты классификации. В поле **Независимые переменные** необходимо перенести одну или несколько независимых переменных.

Параметр **Первая переменная принудительно** позволяет задать первую переменную из списка независимых переменных как первую переменную расщепления.

Поле **Переменная влияния** позволяет указать переменную, которая будет определять, насколько большое влияние данное наблюдение оказывает на процесс построения дерева. Наблюдения с более низкими значениями переменной влияния будут иметь меньшее влияние, а наблюдения с более высокими значениями – большее влияние. При этом значения переменной влияния должны быть положительными.

Выпадающий список **Метод построения** позволяет выбрать метод построения дерева.

В правой части окна находятся пять кнопок, используемых для настройки процедуры **Дерева классификации**.

Кнопка **Вывод** задает появление дерева решений и генерацию таблиц. Можно запросить дополнительную статистическую информацию о модели, графическую интерпретацию соответствующих статистик, также можно запросить генерацию правил классификации для модели в SPSS синтаксисе, в SQL или в обычном текстовом формате.

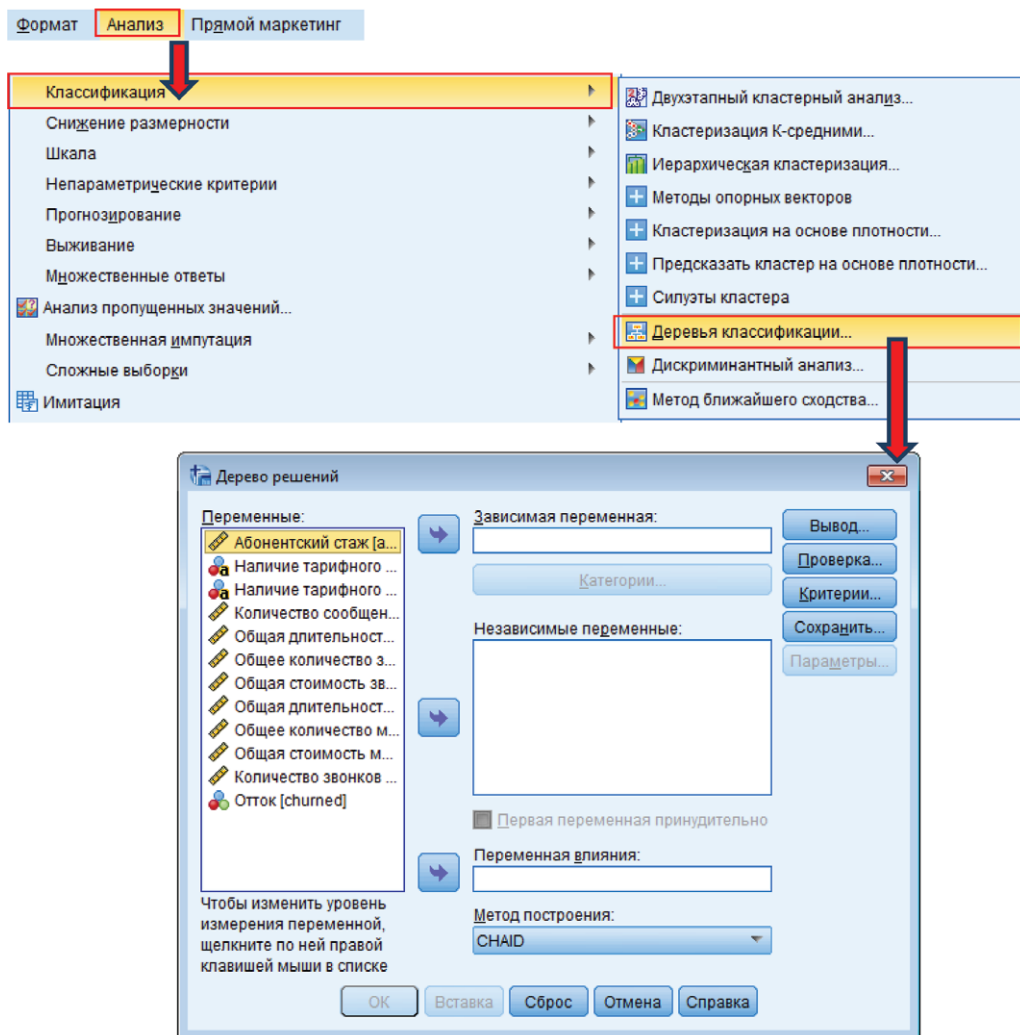


Рис. 2.1 ❖ Запуск процедуры Деревья классификации

Кнопка **Проверка** позволяет построить модель, отобрав только часть данных, и затем посмотреть, как она работает на оставшейся части, которая была исключена при построении модели.

Кнопка **Критерии** задает значения, которые используются в построении модели, такие как минимальное количество наблюдений в каждой группе или сегменте и уровень значимости, используемый в статистических тестах.

Кнопка **Сохранить** добавляет в активный набор данных переменные результатов анализа для каждого наблюдения:

- номер узла, к которому относится наблюдение;
- спрогнозированное значение зависимой переменной (для количественной зависимой переменной сохраняется спрогнозированное среднее значение, для категориальной зависимой переменной – спрогнозированная категория);

- спрогнозированные вероятности категорий зависимой переменной (только для категориальной зависимой переменной);
- принадлежность к обучающей или контрольной выборке.

Кнопка **Параметры** позволяет задать стоимости ошибочной классификации, априорные вероятности, прибыль и затраты по результатам классификации.

## 2.2. Четыре метода деревьев решений

Выпадающий список **Метод построения** в диалоговом окне **Дерево решений** (рис. 2.2) позволяет вам выбрать четыре метода деревьев решений: **CHAID** (используется по умолчанию), **Исчерпывающий CHAID**, **CRT**, **QUEST**.

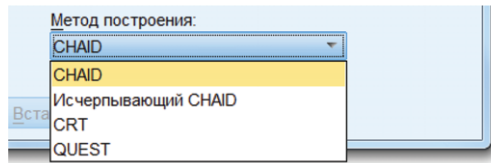


Рис. 2.2 ❖ Выпадающий список **Метод построения**

**CHAID** (расшифровывается как *Chi-square Automatic Interaction Detector – Автоматический обнаружитель взаимодействий*) используется процедурой **Дерева классификации** по умолчанию. Он был разработан Гордоном Каасом в 1980 году и представляет собой метод на основе дерева решений, который исследует взаимосвязь между предикторами и зависимой переменной с помощью статистических тестов.

Каждый раз для разбиения узла выбирается предиктор, сильнее всего взаимодействующий с зависимой переменной. При этом категории каждого предиктора объединяются, если они не имеют между собой статистически значимых отличий по отношению к зависимой переменной, остальные категории рассматриваются как отдельные. Для количественной зависимой переменной используется F-тест, для категориальной зависимой переменной – хи-квадрат Пирсона или хи-квадрат отношения правдоподобия.

Зависимая переменная и предикторы могут быть измерены в номинальной, порядковой и количественной шкалах, при этом количественные предикторы преобразовываются в порядковые переменные. CHAID позволяет осуществлять многомерные расщепления узлов. Каждый узел при разбиении может иметь более 2 потомков, поэтому CHAID имеет тенденцию выращивать более раскидистые деревья, чем бинарные методы. Вместе с тем из-за жестких статистических критериев расщепления нередко дерево CHAID получается нереалистично коротким и тривиальным («грубое» дерево), поэтому требуется тонкая настройка уровней значимости для объединения категорий и разбиения узлов. По сравнению с другими методами, CHAID характеризуется умеренным временем вычислений.

Помимо прочего, метод CHAID обладает собственным способом обработки пропущенных значений. Пропуски рассматриваются как отдельная фактическая категория. В ряде случаев это имеет смысл. Например, отказ отвечать на вопрос о доходе или занятости может оказаться предсказательной категорией для зависимой переменной.

**Исчерпывающий CHAID** является модификацией метода CHAID, предложенной Дэвидом Биггсом, Барри Де Виллем и Эдом Суеном в 1991 году. Он был разработан для устранения недостатка CHAID – ограниченного набора расщеплений для предиктора.

CHAID прекращает объединение категорий, когда обнаруживает, что все оставшиеся категории статистически различаются между собой. Исчерпывающий CHAID исправляет это, продолжая объединять категории предиктора до тех пор, пока не останутся только две суперкатегории. Таким образом, он позволяет найти наилучшее расщепление для каждого предиктора и затем выбрать, какой предиктор нужно расщепить.

Исчерпывающий CHAID идентичен CHAID с точки зрения используемых зависимой переменной и предикторов, статистических тестов значимости взаимодействия и способа обработки пропущенных значений. Вместе с тем, поскольку объединение категорий осуществляется более тщательно, чем в методе CHAID, исчерпывающий CHAID требует большего времени вычислений. Надежность результатов исчерпывающего CHAID выше, чем у CHAID.

**CRT** (расшифровывается как Classification and Regression Tree – Деревья классификации и регрессии) был разработан в 1974–1984 годах профессорами статистики Лео Брейманом (Калифорнийский университет в Беркли), Джеромом Фридманом (Стэнфордский университет), Ричардом Олшеном (Калифорнийский университет в Беркли) и Чарльзом Стоуном (Стэнфордский университет).

Для построения дерева метод CRT использует принцип уменьшения неоднородности в узле. Расщепление узла происходит так, чтобы узел-потомок был более однородным, чем его узел-родитель. В абсолютно однородном узле все наблюдения имеют одно и то же значение целевой переменной (все объекты принадлежат к одной и той же категории целевой переменной). Такой узел еще называют «чистым».

Зависимая переменная может быть измерена в номинальной, порядковой и количественной шкалах. Предикторы могут быть измерены в номинальной, порядковой и количественной шкалах (подробнее о типах шкал читайте в разделе 2.3 «*Шкалы переменных*»). CRT позволяет только одномерные расщепления узлов. Каждый узел при разбиении может иметь лишь 2 потомков. Поэтому CRT имеет тенденцию выращивать высокие деревья с большим количеством уровней. Часто деревья CRT получаются слишком детализированными, имеют много узлов и ветвей, сложны для интерпретации, при этом усложнение дерева не приводит к повышению прогностической способности дерева. Для упрощения структуры дерева и повышения качества модели в методе CRT предусмотрена возможность отсека ветвей (прунинг). Прунинг позволяет получить дерево «подходящего размера», избежать построения ветвистых, усложненных деревьев и при этом достичь лучшего качества модели.

Для обработки наблюдений, у которых пропущено значение в предикторе, используются суррогаты – другие предикторы, имеющие сильную корреляцию с исходной независимой переменной. Таким образом, разбиение, задаваемое суррогатом, будет наиболее близко к разбиению, задаваемому исходным предиктором, по которому имеются пропуски. Метод CRT требует больше время вычислений, по сравнению с другими методами.

**QUEST** (расшифровывается как *Quick, Unbiased, Efficient Statistical Tree* – Быстрое, несмещенное, эффективное статистическое дерево) был предложен в 1997 году профессорами статистики Вэй Ин Ло (Университет Висконсина-Мэдисона) и Ю Шан Ши (Национальный университет Чун Чен, Тайвань).

Метод **QUEST** строит дерево следующим образом: для отбора предикторов используются статистические тесты значимости взаимодействия между зависимой переменной и предиктором, а разбиение узлов задается путем выполнения квадратичного дискриминантного анализа с использованием отобранного предиктора. Зависимая переменная может быть измерена только в номинальной шкале. Предикторы могут быть измерены в номинальной, порядковой и количественной шкалах.

**QUEST** имеет схожие с **CRT** характеристики:

- позволяет только одномерные расщепления узлов;
- каждый узел при разбиении может иметь лишь 2 потомков;
- есть возможность отсечения ветвей (прунинг);
- для обработки наблюдений, у которых пропущено значение в предикторе, используются суррогаты – другие предикторы, имеющие сильную корреляцию с исходной независимой переменной.

Ниже на рис. 2.3 приводится таблица сходств и различий между четырьмя методами деревьев решений, предлагаемых процедурой **Деревья классификации**.

Характеристика метода	CHAID	Exhaustive CHAID	CRT	QUEST
Категориальная зависимая переменная	Да	Да	Да	Да, только номинальная
Категориальные предикторы	Да	Да	Да	Да
Количественная зависимая переменная	Да	Да	Да	Нет
Количественные предикторы	Да, преобразуются в порядковые	Да, преобразуются в порядковые	Да	Да
Тип разбиения	Множественный	Множественный	Бинарный	Бинарный
Цены ошибочной классификации (Построение дерева)	Нет	Нет	Да	Да
Статистические тесты (Отбор предикторов)	Да	Да	Нет	Да
Статистические тесты (Разбиение)	Да	Да	Нет	Нет
Время вычислений	Умеренное	Умеренное	Большое	Умеренное/Большое
Использование априорных вероятностей	Нет	Нет	Да	Да
Пропущенные значения в предикторах	Да, как категория	Да, как категория	Нет, для разбиения используется заместитель	Нет, для разбиения используется заместитель

Рис. 2.3 ❖ Четыре метода деревьев решений

## 2.3. Шкалы переменных

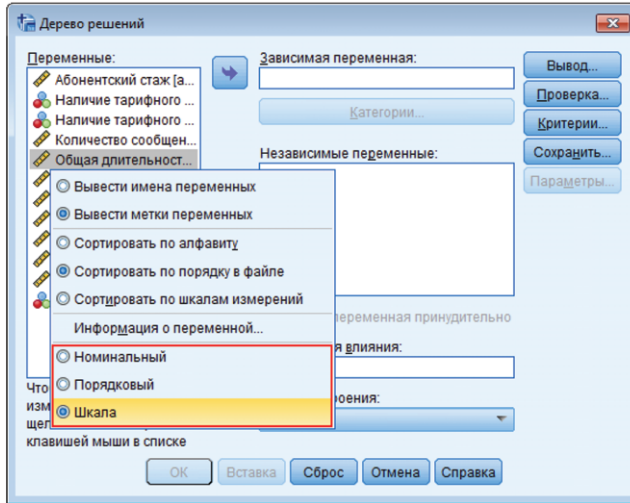
В зависимости от шкалы (уровня измерения) зависимой переменной и независимых переменных деревья решений применяют различные критерии для отбора предикторов и разбиения узлов. Поэтому важно задать правильную шкалу переменной. В IBM SPSS Statistics существуют три типа шкалы: количественная, порядковая, номинальная.

Количественная шкала содержит информацию о расстояниях между уровнями переменной, порядке уровней и количестве объектов в уровнях. Пример предиктора с количественной шкалой – переменная *Возраст*. Например, мы знаем, что расстояние между 25 и 30 в два раза меньше, чем расстояние между 30 и 40, 30-летний на 5 лет старше 25-летнего. Мы можем упорядочить уровни по нарастанию или убыванию интенсивности определенного признака (например, по увеличению возраста): после 25 следует 30, и 30-летний старше 25-летнего. Наконец, мы можем сказать, сколько в выборке человек с тем или иным уровнем (значением) возраста.

Порядковая шкала содержит информацию о порядке уровней и количестве объектов в уровнях. Пример предиктора с порядковой шкалой – переменная *Доход*, разбитая на уровни *низкий*, *средний*, *высокий*. Здесь уже нельзя сказать, что расстояние между уровнями *низкий* и *средний* больше или меньше в определенное количество раз расстояния между уровнями *средний* и *высокий*. Мы не можем утверждать, что человек со средним доходом на  $n$ -ное количество единиц богаче, чем человек с низким доходом. Однако можно упорядочить уровни по нарастанию или убыванию интенсивности определенного признака: сначала следует уровень *низкий*, затем уровень *средний*, и потом уровень *высокий*. Респонденты, относящиеся к уровню *средний*, обладают меньшим доходом, по сравнению с респондентами, относящимися к уровню *высокий*, то есть демонстрируют меньшую интенсивность признака. Также мы можем сказать, сколько в выборке человек с тем или иным уровнем дохода.

Номинальная шкала содержит только информацию о количестве объектов в уровнях. Пример предиктора с номинальной шкалой – переменная *Регион*, который имеет уровни *Алтайский край*, *Новосибирская область*, *Красноярский край*, *Кемеровская область*. Мы ничего не можем сказать о расстояниях между уровнями, о порядке уровней. Мы можем лишь судить о количестве респондентов, проживающих в каждом регионе.

Для изменения шкалы непосредственно в диалоговом окне **Дерево решений** щелкните правой кнопкой мыши по переменной, уровень измерения которой вы хотите изменить. В появившемся контекстном меню, изображенном на рис. 2.4, можно объявить переменную как номинальную, порядковую или количественную. Обратите внимание, что вы можете менять шкалу переменной только до ее перемещения в область **Зависимая переменная** или в область **Независимые переменные**.



**Рис. 2.4** ❖ Изменение шкалы переменной в диалоговом окне **Дерево решений**

Также обратите внимание, что пиктограммы, сопровождающие каждую переменную в диалоговом окне, показывают текущую шкалу этой переменной. Пиктограммы и соответствующие им шкалы показаны ниже на рис. 2.5.



**Рис. 2.5** ❖ Три пиктограммы переменных (слева направо: количественная, номинальная, порядковая)

Кроме того, задать шкалу можно в Редакторе переменных (рис. 2.6). Для этого перейдите в Редактор переменных, щелкнув в левом нижнем углу Редактора данных вкладку **Представление Переменные**. Затем щелкните левой клавишей мыши по полю **Шкала** напротив интересующей переменной и выберите нужную шкалу.



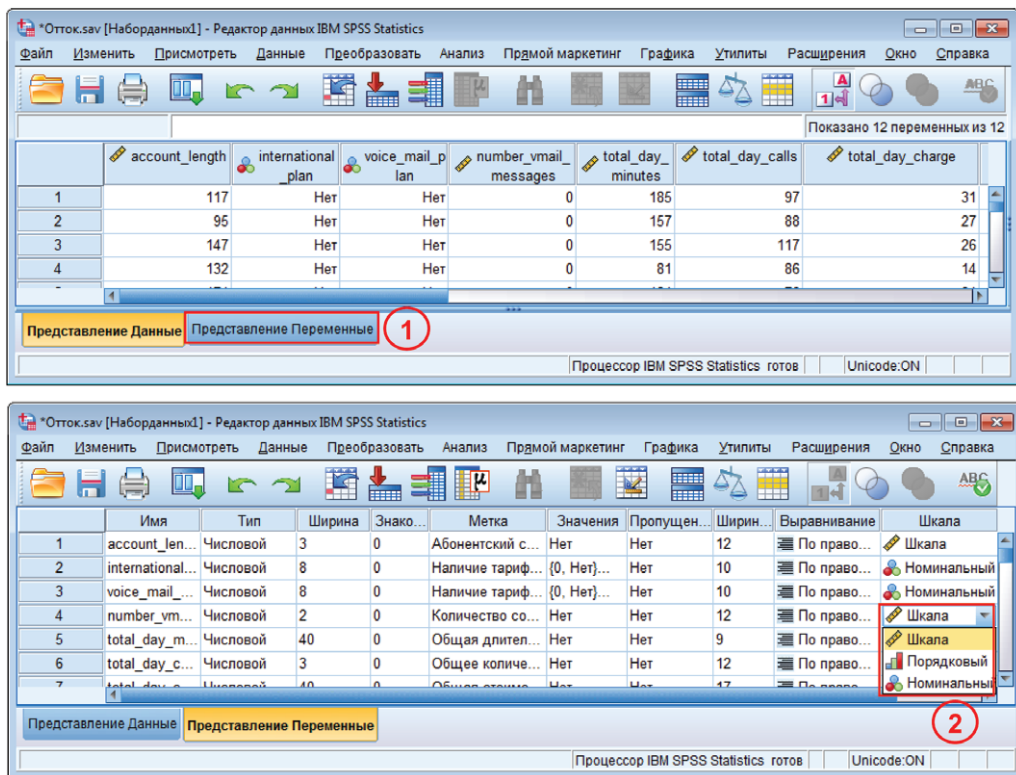


Рис. 2.6 ❖ Изменение шкалы переменной

## 2.4. Определение необходимого размера выборки

Формирование выборки в ходе разработки прогнозной модели обусловлено двумя задачами, а именно разбиением исторической выборки на обучающую и контрольную и распределением наблюдений в категориях зависимой переменной.

Предположим, мы хотим быть уверенным на 95%, что соотношение «плохих» и «хороших» заемщиков в выборке должно отражать генеральную популяцию заемщиков. Зная долю интересующего признака (например, долю «плохих») в популяции, задав z-статистику, определяемую в зависимости от уровня значимости, и ширину доверительного интервала (точность оценки), вы можете найти минимальный объем выборки по нижеприведенной формуле:

$$n = (z^1 p(1 - p))/d^2,$$

где  $z$  – z-статистика для уровня значимости (например, 1,96 для уровня значимости 0,05);  $p$  – доля признака в популяции (например, доля «плохих»);  $d$  – половина ширины доверительного интервала для доли признака (точность оценки доли признака).

Например, предположив наихудший сценарий 50/50 (задающий максимально возможный размер выборки) и точность оценки 5%, или 0,05, вычислим минимальный объем выборки:

$$n = (1,962 \times 0,5 \times (1 - 0,5))/0,05^2 = 385.$$

Однако нужно помнить, что эти методы дают минимальные размеры выборки. Необходимо придерживаться такого метода формирования выборки, который считается наиболее обоснованным до тех пор, пока размеры выборки позволяют убедиться в удовлетворительных статистических и практических результатах.

На практике обычно используют правило 20 EPV (Event Per Variable), сформулированное американским статистиком Фрэнком Харреллом<sup>1</sup>. Оно связывает минимальный объем выборки с количеством наблюдений в миноритарной (наименьшей по размеру) категории зависимой переменной и количеством предикторов, поданных на вход модели.

Согласно этому правилу, необходимо взять количество наблюдений в исторической выборке, относящихся к миноритарной категории зависимой переменной (в кредитном скоринге это «плохие» заемщики). Это число наблюдений нужно разделить на количество заданных предикторов. На один предиктор должно приходиться не менее 20 наблюдений. Если это правило выполняется, то объем выборки достаточный.

Например, у нас есть выборка в 2000 наблюдений (1700 «хороших» и 300 «плохих») и 16 независимых переменных. Проверяя выполнение правила 20EPV, мы получаем  $300/16 = 18,75$ . Наша выборка не обеспечивает достаточного количества наблюдений в миноритарной категории зависимой переменной.

Если планируется проверка модели с разбиением набора данных на обучающую и контрольную выборки (например, 70%:30%), необходимо руководствоваться еще более жестким правилом: взять число наблюдений в миноритарной категории зависимой переменной по контрольной выборке и разделить на количество заданных предикторов.

Например, имеется историческая выборка объемом 5000 наблюдений (4000 «хороших» и 1000 «плохих») и 10 независимых переменных. Теперь разбиваем историческую выборку случайным образом на обучающую и контрольную. В обучающую выборку попали 3500 наблюдений (2800 «хороших» и 700 «плохих»), а в контрольную выборку попали 1500 наблюдений (1200 «хороших» и «300 плохих»). Проверяя выполнение правила 20EPV, мы получаем  $300/10 = 30$ . Историческую выборку данного объема можно использовать для моделирования.

## 2.5. Знакомство с методом CHAID

### 2.5.1. Описание алгоритма

Перед началом работы алгоритм CHAID преобразует все имеющиеся количественные предикторы в порядковые переменные, имеющие по умолчанию 10 приблизительно равных по объему категорий (число категорий можно устанавливать самостоятельно).

Затем алгоритм приступает к построению дерева, итеративно применяя к каждому узлу, начиная с корневого, процедуры объединения категорий, расщепления узла и проверки правил остановки.

<sup>1</sup> <http://biostat.mc.vanderbilt.edu/wiki/Main/ManuscriptChecklist>.

### Объединение категорий

1. Для каждого предиктора с числом категорий больше двух<sup>1</sup> алгоритм ищет пару категорий (для порядковых переменных можно брать лишь две смежные категории, для номинальных переменных – любые две категории), меньше всего статистически различающихся по зависимой переменной. Для этого он выполняет статистические тесты. Выбор теста определяется типом шкалы зависимой переменной. Для номинальной зависимой переменной используется тест хи-квадрат Пирсона, или хи-квадрат отношения правдоподобия. Алгоритм строит двухвходовую таблицу сопряженности с категориями предиктора в качестве строк и категориями зависимой переменной в качестве столбцов. Он проверяет нулевую гипотезу о том, что категории предиктора не отличаются друг от друга с точки зрения распределения категорий зависимой переменной. Для количественной зависимой переменной используется F-тест. Алгоритм осуществляет дисперсионный анализ и проверяет нулевую гипотезу о том, что средние значения зависимой переменной для различных категорий предиктора не различаются между собой. Для порядковой зависимой переменной используется тест хи-квадрат отношения правдоподобия. Алгоритм подгоняет модель эффектов строк, где строки являются категориями предиктора, а столбцы – категориями зависимой переменной. В рамках каждого теста алгоритм вычисляет  $p$ -значение – вероятность того, что случайная величина с распределением тестовой статистики при нулевой гипотезе примет значение, не меньшее, чем фактическое значение тестовой статистики. Таким образом, задача алгоритма сводится к тому, чтобы для каждого предиктора найти пару категорий с наибольшим  $p$ -значением, поскольку именно такие категории будут меньше всего статистически различаться по зависимой переменной.

2. Найдя наибольшее  $p$ -значение для пары категорий, алгоритм сравнивает его с заданным уровнем значимости для объединения категорий.

Если  $p$ -значение:

- меньше или равно заданному уровню значимости для объединения категорий – алгоритм переходит к вычислению скорректированных  $p$ -значений для полученного набора категорий (шаг 3);
- больше уровня значимости для объединения категорий – эта пара объединяется в отдельную составную категорию, в результате формируется новый набор категорий предиктора и процесс начинается заново с поиска пары категорий с наибольшим  $p$ -значением.

### ПРИМЕЧАНИЕ

Чтобы задать уровень значимости для объединения категорий, необходимо выполнить следующие действия:

- в главном диалоговом окне **Дерево решений** нажмите кнопку **Критерии**;
- в открывшемся диалоговом окне **Деревья решений: Критерии** откройте вкладку **CHAID**;
- в панели **Уровни значимости для** выберите для параметра **Объединения категорий** необходимое значение (значение по умолчанию 0,05).

Более подробно об этом читайте в разделе 3.1 «**Построение деревьев CHAID с измененными критериями**».

<sup>1</sup> Если предиктор имеет одну категорию, он исключается из анализа. Если предиктор имеет две категории, переходит к шагу 3.

(Опционно) Если новая составная категория состоит из трех и более исходных категорий, алгоритм находит внутри этой составной категории наилучшее бинарное расщепление, которое дает наименьшее  $p$ -значение. Алгоритм выполняет бинарное расщепление, если его  $p$ -значение не превышает уровня значимости для разбиения объединенных категорий.

### ПРИМЕЧАНИЕ

Чтобы задать разбиение объединенных категорий, необходимо выполнить следующие действия:

- в главном диалоговом окне **Дерево решений** нажмите кнопку **Критерии**;
- в открывшемся диалоговом окне **Деревья решений: Критерии** откройте вкладку **CHAID**;
- отметьте параметр **Допускать разбиение объединенных категорий в узле** (по умолчанию не используется).

Более подробно об этом читайте в разделе 3.1 «**Построение деревьев CHAID с измененными критериями**».

3. Для сформированного набора категорий предиктора алгоритм вычисляет скорректированное  $p$ -значение как  $p$ -значение, умноженное на поправку Бонферрони. Поправка Бонферрони представляет собой число возможных способов, с помощью которых исходные категории предиктора могут быть объединены в итоговые категории.

### ПРИМЕЧАНИЕ

Чтобы отменить применение поправки Бонферрони, необходимо выполнить следующие действия:

- в главном диалоговом окне **Дерево решений** нажмите кнопку **Критерии**;
- в открывшемся диалоговом окне **Деревья решений: Критерии** откройте вкладку **CHAID**;
- деактивируйте параметр **Корректировать уровни значимости с использованием поправки Бонферрони** (по умолчанию активирован).

### *Расщепление узла*

После вычисления скорректированных  $p$ -значений для итоговых наборов категорий по всем предикторам алгоритм переходит к этапу расщепления узла.

1. На этапе расщепления алгоритм выбирает, какой предиктор обеспечит наилучшее разбиение узла. Для этого предиктор должен иметь наименьшее скорректированное  $p$ -значение (то есть является наиболее статистически значимым).

2. Найдя предиктор с наименьшим скорректированным  $p$ -значением, алгоритм сравнивает его с заданным уровнем значимости для расщепления.

Если  $p$ -значение:

- меньше или равно заданному уровню значимости для расщепления – алгоритм разбивает узел с использованием данного предиктора (категории предиктора становятся дочерними узлами);
- больше заданного уровня значимости для расщепления, то алгоритм не расщепляет узел и узел рассматривается как терминальный.