

Содержание

Предисловие	7
Глава 1. Введение в метод деревьев решений	9
1.1. Введение в методологию деревьев решений.....	9
1.2. Преимущества и недостатки деревьев решений.....	11
1.3. Задачи, выполняемые с помощью деревьев решений	12
Вопросы к главе 1.....	14
Часть I. ПОСТРОЕНИЕ ДЕРЕВЬЕВ РЕШЕНИЙ В IBM SPSS STATISTICS	16
Глава 2. Основы прогнозного моделирования с помощью деревьев решений CHAID	17
2.1. Запуск процедуры Деревья классификации.....	17
2.2. Четыре метода деревьев решений.....	19
2.3. Шкалы переменных.....	21
2.4. Определение необходимого размера выборки.....	23
2.5. Знакомство с методом CHAID	25
2.5.1. Описание алгоритма	25
2.5.2. Немного о тесте хи-квадрат.....	28
2.5.3. Немного об F-тесте	28
2.5.4. Способы объединения категорий предикторов.....	29
2.5.5. Поправка Бонферрони.....	30
2.5.6. Иллюстрация работы CHAID на конкретном примере	30
2.6. Построение дерева классификации CHAID.....	35
2.6.1. Настройка процедуры Деревья классификации.....	35
2.6.2. Работа с отчетом о построении модели	36
2.7. Работа с прогнозами модели	43
2.7.1. Получение результатов классификации	43
2.7.2. Сохранение прогнозов модели в файле данных.....	44
2.7.3. Самостоятельное построение таблицы классификации и изменение порогового значения вероятности	48
2.8. Анализ ROC-кривой	57
2.8.1. Терминология анализа ROC-кривой	57
2.8.2. Оценка дискриминирующей способности модели и выбор порогового значения с помощью ROC-кривой	61
2.9. Проверка модели	67
2.9.1. Методы проверки модели в процедуре Деревья классификации.....	67
2.9.2. Работа с результатами проверки модели	70
2.10. Дополнительные настройки вывода результатов	84
2.10.1. Настройки вывода дерева	84
2.10.2. Построение таблицы дерева.....	85

2.10.3. Настройки вывода статистик	86
2.10.4. Построение таблиц выигрышей для узлов и перцентилей.....	88
2.10.5. Настройки вывода графиков	89
2.10.6. Построение графиков выигрышей, индексов и откликов	91
2.10.7. Настройки вывода правил классификации.....	93
2.10.8. Применение правил классификации к новому набору данных.....	95
2.11. Построение дерева регрессии CHAID.....	104
2.12. Использование принудительной переменной расщепления	108
Вопросы к главе 2.....	109

Глава 3. Продвинутое моделирование с помощью деревьев решений CHAID 112

3.1. Построение деревьев CHAID с измененными критериями.....	112
3.1.1. Настройка правил остановки.....	112
3.1.2. Построение деревьев CHAID с измененными правилами остановки.....	113
3.1.3. Настройка статистических тестов для разбиения узлов и объединения категорий предикторов.....	119
3.1.4. Построение дерева CHAID с измененными статистическими тестами	120
3.1.5. Настройка обработки количественных предикторов.....	121
3.1.6. Построение дерева CHAID с измененным числом интервалов для количественных предикторов	122
3.2. Метод Исчерпывающий CHAID	123
3.3. Обзор параметров деревьев решений.....	124
3.4. Работа с пропусками в методе CHAID	126
3.4.1. Настройка обработки пропущенных значений	126
3.4.2. Построение дерева CHAID на основе данных, содержащих пропуски	129
3.5. Работа со стоимостями ошибочной классификации в методе CHAID.....	130
3.5.1. Настройка стоимостей ошибочной классификации.....	130
3.5.2. Построение дерева CHAID с измененными стоимостями ошибочной классификации	133
3.6. Работа с прибылями в методе CHAID.....	136
3.6.1. Настройка прибылей	136
3.6.2. Построение дерева CHAID с заданными значениями прибыли	137
3.7. Работа со значениями	141
3.8. Применение метода CHAID для биннинга переменных (на примере конкурсной задачи ОТП Банка).....	144
3.8.1. Преимущества и недостатки биннинга	144
3.8.2. Предварительная подготовка данных.....	146
3.8.3. Определение важности переменных с помощью случайного леса.....	162
3.8.4. Анализ мультиколлинеарности.....	165
3.8.5. Выполнение автоматического биннинга переменных	167
3.8.6. Построение моделей логистической регрессии на основе исходных предикторов и предикторов, категоризированных с помощью CHAID.....	169

3.8.7. Выполнение биннинга переменных с помощью процедуры Оптимальная категоризация.....	172
3.8.8. Построение модели логистической регрессии на основе оптимально категоризированных предикторов	175
3.8.9. Преобразование количественных переменных для максимизации нормальности	176
3.8.10. Построение модели логистической регрессии с использованием CHAID и преобразования корня третьей степени	180
3.9. Построение ансамбля логистической регрессии и дерева CHAID.....	181
Вопросы к главе 3.....	186

Глава 4. Построение деревьев решений CRT и QUEST 188

4.1. Знакомство с методом CRT	188
4.1.1. Описание алгоритма	189
4.1.2. Неоднородность.....	190
4.1.3. Внутриузловая дисперсия	191
4.1.4. Метод отсечения ветвей на основе меры стоимости-сложности	192
4.1.5. Обработка пропущенных значений.....	193
4.1.6. Иллюстрация работы CRT на конкретном примере	193
4.2. Построение дерева классификации CRT.....	195
4.3. Построение дерева CRT с измененными критериями	199
4.3.1. Настройка мер неоднородности для отбора предикторов и расщепления узлов	199
4.3.2. Настройка отсечения ветвей.....	200
4.3.3. Построение дерева CRT с последующим отсечением ветвей.....	201
4.3.4. Настройка суррогатов для обработки пропущенных значений	203
4.3.5. Построение дерева CRT на основе данных, содержащих пропуски	203
4.4. Вывод важности предикторов.....	206
4.5. Работа с априорными вероятностями в методе CRT.....	207
4.5.1. Настройка априорных вероятностей	207
4.5.2. Построение дерева CRT с измененными априорными вероятностями.....	208
4.6. Знакомство с методом QUEST.....	210
4.6.1. Описание алгоритма	211
4.6.2. Метод отсечения ветвей на основе меры стоимости-сложности	213
4.7. Построение дерева классификации QUEST	213
4.8. Сравнение метода QUEST с другими методами деревьев решений	216
4.9. Построение дерева QUEST с измененными критериями.....	216
4.9.1. Настройка статистических тестов для отбора предикторов	217
4.9.2. Построение дерева QUEST с последующим отсечением ветвей	217
Вопросы к главе 4.....	219

Глава 5. Редактор дерева 220

5.1. Просмотр диаграммы дерева в Редакторе	220
5.2. Просмотр содержимого узла в Редакторе.....	221
5.3. Настройка внешнего вида диаграммы дерева в Редакторе.....	222

5.4. Изменение ориентации диаграммы дерева в Редакторе.....	223
5.5. Настройка содержимого узла в Редакторе.....	223
5.6. Отбор наблюдений в Редакторе.....	224
5.7. Иллюстрация работы в Редакторе дерева на конкретном примере.....	225

Часть II. ПОСТРОЕНИЕ ДЕРЕВЬЕВ РЕШЕНИЙ В R 229

Глава 6. Построение деревьев решений CHAID с помощью пакета R CHAID 230

6.1. Построение и интерпретация дерева классификации CHAID.....	230
6.2. Работа с прогнозами модели.....	234
6.3. Сохранение результатов прогноза.....	239
6.4. Применение модели к новым данным.....	239
6.5. Проверка модели.....	241
6.6. Биннинг переменных.....	242
6.6.1. Биннинг в пакете rattle.....	242
6.6.2. Биннинг в пакете smbinning.....	244
Вопросы к главе 6.....	252

Глава 7. Построение деревьев решений CRT с помощью пакета R rpart..... 253

7.1. Метод отсечения ветвей на основе стоимости-сложности с кросс-проверкой.....	253
7.2. Построение и интерпретация дерева классификации CRT.....	254
7.3. Прунинг дерева CRT.....	262
7.4. Работа с прогнозами модели.....	264
7.5. Сохранение результатов прогноза.....	267
7.6. Применение модели к новым данным.....	268
7.7. Построение и интерпретация дерева регрессии CRT.....	269
Вопросы к главе 7.....	273

Ключи к вопросам..... 275

Предисловие

Данная книга открывает серию пособий, посвященных практическому применению методов машинного обучения на базе популярных статистических пакетов IBM SPSS Statistics и R. В первом выпуске освещается метод деревьев решений. Деревья решений – это эффективный метод машинного обучения, использующийся в прогнозном моделировании. Кроме того, при решении задач бинарной классификации он нередко дополняет метод логистической регрессии. Аналитики кредитного бюро TransUnion для построения скоринговых моделей используют логистическую регрессию, а для отбора переменных в модель логистической регрессии (рассматриваются сотни переменных) – деревья решений CRT или случайный лес. Наша компания при построении прогнозных моделей на основе логистической регрессии использует метод деревьев решений, чтобы сформировать новые переменные для лучшего прогнозирования дефолта. Аналитики Citibank USA разбивают популяцию заемщиков на сегменты, применяя дерево решений, а затем в каждом сегменте строят модели доходности с помощью линейной регрессии или модели риска с помощью логистической регрессии.

Кратко о самой книге. В ней я детально расскажу о том, как строить деревья решений, интерпретировать их, оценивать дискриминирующую способность полученных моделей, улучшать их, сохранять результаты и применять правила классификации/прогноза, полученные с помощью дерева, к новым данным. Кроме того, я расскажу о том, как с помощью дерева решений улучшить модель логистической регрессии. Глава 1 кратко знакомит с терминологией метода деревьев решений, в ней рассказывается о преимуществах и недостатках деревьев, задачах, которые можно выполнить с их помощью. Главы 2–4 посвящены построению деревьев решений в IBM SPSS Statistics 24.0. В главе 2 освещается CHAID – один из самых популярных методов деревьев решений. В главе 3 я покажу, как можно менять параметры дерева CHAID, влияя на результаты классификации. Здесь же я расскажу о том, как можно выполнить биннинг переменных для включения в модель логистической регрессии, используя дерево решений CHAID и случайный лес. Для иллюстрации выбрана конкурсная задача предсказания отклика ОТП Банка. Кроме того, на данном примере я покажу, как выполняется предварительная подготовка данных и решаются вопросы, связанные с автоматизацией построения моделей (для этого будет использован командный синтаксис SPSS). Код, автоматизирующий процесс построения прогнозных моделей, вы можете в дальнейшем использовать в собственных проектах. В этой же главе будет рассмотрена разработка ансамбля модели логистической регрессии и дерева CHAID. Глава 4 посвящена методам деревьев CRT и QUEST. В главе 5 рассказывается о Редакторе дерева. В главах 6 и 7 я подробно рассмотрю процесс построения и интерпретации деревьев решений в пакетах R CHAID и gpart. Всю необходимую информацию об IBM SPSS Statistics вы найдете на официальном сайте компании IBM. Информацию о программном пакете R можно найти на официальной странице проекта R. Все вопросы, возникшие в ходе чтения книги, можно направлять по адресу info@gewissta.ru.

Освещаемые темы будут интересны маркетологам, риск-аналитикам и другим специалистам, занимающимся разработкой и внедрением прогнозных моделей.

В заключение я хочу поблагодарить моих взыскательных клиентов и коллег из TransUnion, DBS Bank и Citibank USA, в особенности Дмитрия Майорова (Citibank N. A., ArrowModel) и Барри Уилка (Google) за их ценные советы и замечания, высказанные в ходе подготовки книги.

Артем Груздев,
генеральный директор ИЦ «Гевисста»

Введение в метод деревьев решений

1.1. Введение в методологию деревьев решений

Как и регрессионный анализ, деревья решений являются методом изучения статистической взаимосвязи между одной зависимой переменной и несколькими независимыми (предикторными) переменными. Базовое отличие метода деревьев решений от регрессионного анализа заключается в том, что взаимосвязь между значением зависимой переменной и значениями независимых переменных представлена не в виде общего прогнозного уравнения, а в виде древовидной структуры, которую получают с помощью иерархической сегментации данных.

Берется весь обучающий набор данных, называемый **корневым узлом**, и разбивается на два или более **узлов (сегментов)** так, чтобы наблюдения, попавшие в разные узлы, максимально отличались друг от друга по зависимой переменной (например, выделяем два узла с наибольшим и наименьшим процентами «плохих» заемщиков). В роли **правил разбиения**, максимизирующих эти различия, выступают значения независимых переменных (пол, возраст, доход и др.). Качество разбиения оценивается с помощью статистических критериев. Правила и статистики отмечаются на **ветвях** – линиях, которые соединяют разбиваемый узел с узлами, полученными в результате разбиения. Для каждого узла вычисляются **вероятности** в виде **процентных долей** категорий зависимой переменной (если зависимая переменная является категориальной) или средние значения зависимой переменной (если зависимая переменная является количественной). В результате выносится **решение** – спрогнозированная категория зависимой переменной (если зависимая переменная является категориальной) или спрогнозированное среднее значение зависимой переменной (если зависимая переменная является количественной).

Аналогичным образом каждый узел, получившийся в результате разбиения корневого узла, разбивается дальше на узлы, т. е. узлы внутри узла, и т. д. Этот процесс продолжается до тех пор, пока есть возможность разбиения на узлы. Данный процесс сегментации называется **рекурсивным разделением**. Получившаяся иерархическая структура, характеризующая взаимосвязь между значением зависимой переменной и значениями независимых переменных, называется **деревом**.

Иногда для обозначения разбиваемого узла применяется термин **родительский узел**. Новые узлы, получившиеся в результате разбиения, называются **дочерними**

узлами (или узлами-потомками). Когда впоследствии дочерний узел разбивается сам, он становится родительским узлом. Окончательные узлы, которые в дальнейшем не разбиваются, называются **терминальными узлами** дерева. Их еще называют **листьями**, потому что в них рост дерева останавливается. Лист представляет собой наилучшее окончательное решение, выдаваемое деревом. Здесь мы определяем группы клиентов, обладающие желаемыми характеристиками (например, тех, кто погасит кредит или откликнется на наше маркетинговое предложение).

Обратите внимание, если вы прогнозируете вероятность значения категориальной зависимой переменной по соответствующим значениям предикторов, дерево решений называют **деревом классификации** (рис. 1.1). Например, дерево классификации строится для вычисления вероятности дефолта у заемщика (на основе спрогнозированной вероятности мы относим его к «плохому» или «хорошему» заемщику). Если дерево решений используется для того, чтобы спрогнозировать среднее значение количественной зависимой переменной по соответствующим значениям предикторов, его называют **деревом регрессии** (рис. 1.2). Например, дерево регрессии строится, чтобы вычислить средний размер вклада у клиента.

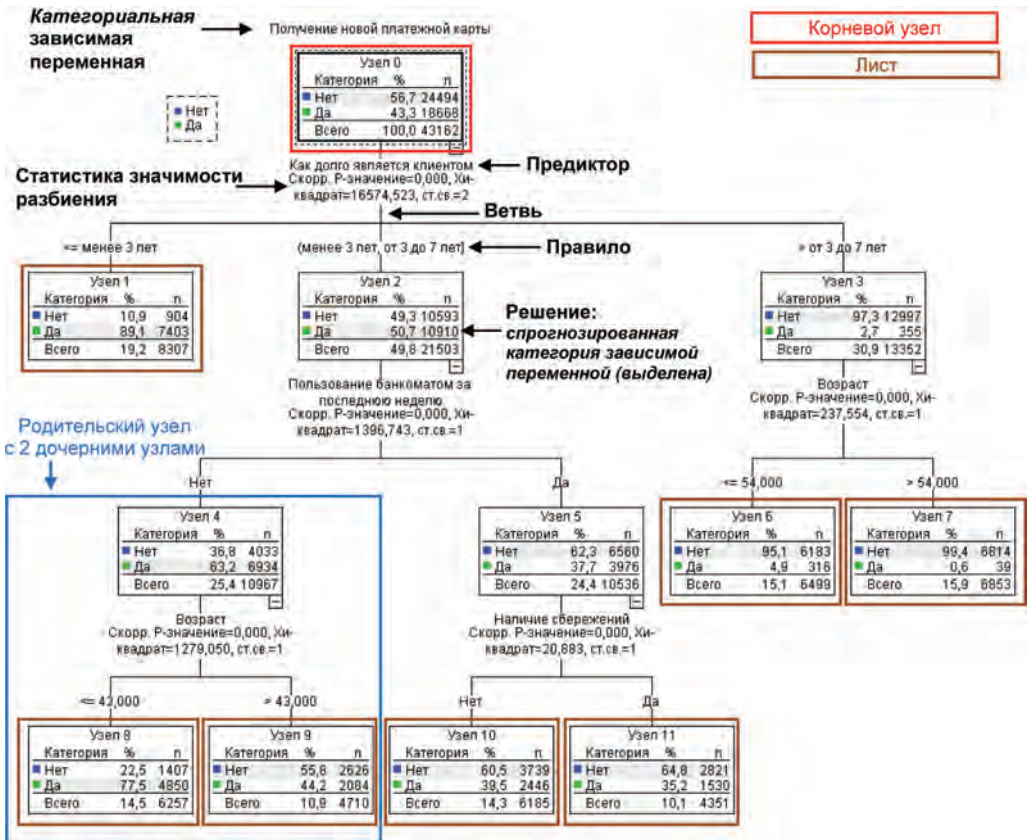


Рис. 1.1 ❖ Пример дерева классификации



Рис. 1.2 ❖ Пример дерева регрессии

1.2. Преимущества и недостатки деревьев решений

Метод деревьев решений обладает рядом преимуществ. Главное из них – это наглядность представления результатов (в виде иерархической структуры дерева). Деревья решений позволяют работать с большим числом независимых переменных. На вход можно подавать все существующие переменные, алгоритм сам выберет наиболее значимые среди них, и только они будут использованы для построения дерева (автоматический отбор предикторов). Деревья решений способны выявлять нелинейные взаимосвязи, сложные взаимодействия, которые нелегко обнаружить в рамках стандартных статистических моделей. Они могут работать с любым типом переменной, таким образом, зависимые и независимые переменные могут быть количественными, порядковыми и номинальными. Деревья решений устойчивы к выбросам, поскольку разбиения основаны на количестве наблюдений внутри диапазонов значений, выбранных для расщепления, а не на абсолютных значениях. Перед построением модели необязательно импутировать пропущенные значения, поскольку деревья используют собственные процедуры обработки пропусков. Требования, выдвигаемые методом деревьев решений к распределению переменных, не являются строгими.

К недостаткам метода деревьев решений можно отнести отсутствие простого общего прогнозного уравнения, выражающего модель (в отличие от регрессионного анализа). Другой недостаток заключается в том, что некоторым методам деревьев решений (например, CRT) свойственно переобучение. Речь идет о ситуации, когда деревья получаются слишком детализированными, имеют много узлов и ветвей, сложны для интерпретации, что требует специальной процедуры отсеечения ветвей (она называется прунинг). Наконец, для методов одиночных деревьев характерна проблема множественных сравнений. Перед расщеплением узла дерево сравнивает различные варианты разбиения, число этих вариантов зависит от числа уровней предикторов, как правило, происходит смещение выбора в пользу переменных, у которых большее ко-

личество уровней. Все это обуславливает определенную нестабильность результатов. Небольшие изменения в наборе данных могут приводить к построению совершенно другого дерева. В силу иерархичности дерева изменения в верхних узлах ведут к изменениям во всех узлах, расположенных ниже. Отмечу, что в большей степени вышесказанное относится к методу CRT. Чтобы достигнуть удовлетворительной прогностической способности CRT, один из его разработчиков – Лео Брейман – пришел к идее случайного леса, когда из обучающего набора извлекаются случайные выборки (того же объема, что и исходный обучающий набор) с возвращением, по каждой строится дерево с использованием случайно отобранных предикторов, и затем результаты, полученные по каждому дереву, усредняются. Однако при таком подходе теряется главное преимущество деревьев решений – простота интерпретации.

1.3. Задачи, выполняемые с помощью деревьев решений

Прежде всего деревья решений используются в маркетинге для сегментации клиентской базы. Например, деревья позволяют определить, какие демографические группы имеют максимальный показатель отклика. Эту информацию можно использовать, чтобы максимизировать отклик при будущей прямой рассылке.

Кроме того, деревья применяются для задач прогнозирования и классификации, когда моделируется взаимосвязь между зависимой переменной и предиктором. С этой точки зрения деревья решений сравнивают с логистической регрессией и линейной регрессией. Деревья решений более эффективны, по сравнению с регрессионным анализом, в тех случаях, когда взаимосвязи между предикторами и зависимой переменной являются нелинейными, переменные имеют несимметричные распределения, наблюдаются большое количество коррелирующих между собой переменных, взаимодействие высоких порядков, аномальные значения. Если же предпосылки регрессионного анализа выполняются, то логистическая регрессия (когда зависимая переменная является категориальной) или линейная регрессия (когда зависимая переменная является количественной) может дать лучший результат. Это обусловлено тем, что деревья пытаются описать линейную связь между переменными путем многократных разбиений по предикторам. CHAID делает это за счет расщепления сразу на несколько категорий, CRT и QUEST пытаются уловить эту связь посредством серии бинарных делений, и это может быть менее эффективно, по сравнению с подбором параметров в регрессионном анализе. Однако проблема заключается в том, что данные обычно содержат как линейные, так и нелинейные зависимости, переменные с симметричными и асимметричными распределениями. Поэтому опытный моделировщик может построить ансамбль логистической регрессии и дерева решений (при условии, что оно использует строгие статистические критерии для отбора предикторов разбиения узлов), чтобы скомпенсировать недостатки обоих методов. Ансамбли дерева решений CHAID и логистической регрессии используются в моделях оттока в телекоме, где данные часто характеризуются переменными с U-образным распределением.

В банковском скоринге деревья решений используются как вспомогательный инструмент при разработке модели логистической регрессии. Приведем конкретные примеры такого применения дерева.

В кредитном скоринге использование нескольких скоринговых карт для одного портфеля обеспечивает лучшее дифференцирование риска, чем использование одной скоринговой карты. Это характерно, когда нам приходится работать с разнородной аудиторией, состоящей из различных групп, и одна и та же скоринговая карта не может работать достаточно эффективно для всех. Например, в скоринге кредитных карточек выделяют сегменты «активные клиенты» и «неактивные клиенты», «клиенты в просрочке» и «клиенты, не имеющие просрочек». Переменные в таких сегментах будут сильно отличаться. Например, для активных кредитных карт утилизация будет сильной переменной, а для неактивных – слабой. И наоборот, может оказаться, что время неактивности для активных клиентов равно 0, а для неактивных клиентов время неактивности окажется сильной переменной. Для этих целей выполняют сегментацию клиентов. Первый способ сегментации – деление на группы на основе опыта и отраслевых знаний с последующей аналитической проверкой. Второй способ – это сегментация с помощью статистических методов типа кластерного анализа или деревьев решений. При этом, по сравнению с кластерным анализом, деревья решений обладают преимуществом: они формулируют четкие правила выделения сегментов, а сами выделенные сегменты статистически значимо отличаются между собой по зависимой переменной. В дальнейшем для каждого из сегментов можно построить собственную модель логистической регрессии, разработать скоринговую карту и сформулировать кредитные правила. В Citibank USA является стандартной практикой делать дерево с двумя-тремя уровнями и в каждом узле подгонять свою модель логистической регрессии. В основе скорингового балла FICO также лежит сегментация на основе деревьев решений. Об эффективности использования сегментации в кредитном скоринге пишет в своей книге «Скоринговые карты для оценки кредитных рисков» известный эксперт по управлению рисками Наим Сиддики¹, а также один из разработчиков алгоритмов скоринга компании FICO Брюс Ходли².

С помощью деревьев решений из большого числа предикторов можно выбрать переменные, полезные для построения модели логистической регрессии. Например, из 100 переменных дерево включило в модель 25 переменных, таким образом, у нас появляется информация о том, какие переменные наверняка можно включить в модель логистической регрессии. Метод деревьев решений CRT позволяет вычислить важность переменных, использованных в модели дерева. Мы уже можем ранжировать переменные по степени полезности.

Деревья решений можно использовать для биннинга – перегруппировки категориального предиктора или дискретизации количественного предиктора с целью лучшего описания взаимосвязи с зависимой переменной. Например, при построении модели логистической регрессии часто обнаруживается, что взаимосвязи между количественным предиктором и интересующим событием являются нелинейными. Уравнение логистической регрессии, несмотря на нелинейное преобразование своего выходного значения (логит-преобразование), все равно моделирует линейные зависимости между предикторами и зависимой переменной. Возьмем пример нелинейной зависимости между стажем работы в банке и внутренним мошенничеством.

¹ *Наим Сиддики*. Скоринговые карты для оценки кредитных рисков. М.: Манн, Иванов и Фабер, 2014.

² *Breiman L.* (2001). Statistical modeling: The two cultures.

Допустим, рассчитанный регрессионный коэффициент в уравнении логистической регрессии получился отрицательным. Это значит, что вероятность совершения внутреннего мошенничества с увеличением стажа работы уменьшается. Однако, выполнив разбивку переменной с помощью дерева CHAID на категории до 12 месяцев, от 12 до 36 месяцев, от 36 до 60 месяцев и больше 60 месяцев, стало видно, что зависимость между стажем и внутренним мошенничеством нелинейная. Первая (до 12 месяцев) и последняя (больше 60 месяцев) категории склонны к внутреннему мошенничеству, а промежуточные сегменты, наоборот, не склонны к внутреннему мошенничеству. После правильной разбивки переменной, проведенной с помощью дерева, связь между предиктором и зависимой переменной становится больше похожа на реальную.

Строя модель логистической регрессии, нередко приходится работать с предикторами, у которых большое количество категорий. Как правило, речь идет о географических переменных (регион, область регистрации, область фактического пребывания заемщика, область торговой точки, где клиент брал кредит) и переменных, фиксирующих профессию или сферу занятости заемщика. Если включить такие переменные в модель логистической регрессии, то переменная с k категориями будет преобразована в $k - 1$ дамми-переменных, которые станут в модели логистической регрессии статистически незначимыми. Только представьте, сколько будет дамми-переменных, если у вас 4 географические переменные с 89 категориями. Исключение таких переменных из анализа также не рационально, поскольку они могут дать ценную информацию. Поэтому можно выполнить биннинг с целью укрупнения категорий, а можно построить по этим четырем переменным дерево решений. В результате дерево укрупнит категории переменных и скомбинирует переменные так, чтобы полученные комбинации характеристик максимизировали различия по зависимой переменной. Таковую переменную, где категориями являются терминальные узлы дерева, можно включить в модель логистической регрессии.

Вопросы к главе 1

1. Терминальный узел – это:
 - а) самый верхний узел, представляющий всю выборку наблюдений;
 - б) узел, в котором рост дерева останавливается;
 - с) любой расщепляемый узел;
 - д) новый узел, появившийся в результате расщепления узла.
2. Деревья классификации строятся:
 - а) для количественной зависимой переменной;
 - б) для категориальной зависимой переменной;
 - с) для порядковой зависимой переменной;
 - д) для номинальной зависимой переменной;
 - е) для бинарной зависимой переменной;
 - ф) для любой зависимой переменной.
3. Деревья регрессии строятся:
 - а) для количественной зависимой переменной;
 - б) для категориальной зависимой переменной;

- с) для номинальной зависимой переменной;
 - д) для бинарной зависимой переменной;
 - е) для любой зависимой переменной.
4. В качестве правил разбиения используются:
- а) значения независимых переменных;
 - б) значения зависимых переменных;
 - с) значения категориальных переменных;
 - д) значения количественных переменных.

ЧАСТЬ |



**ПОСТРОЕНИЕ
ДЕРЕВЬЕВ РЕШЕНИЙ
В IBM SPSS STATISTICS**

Основы прогнозного моделирования с помощью деревьев решений CHAID

2.1. Запуск процедуры Деревья 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... Деревья 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.

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

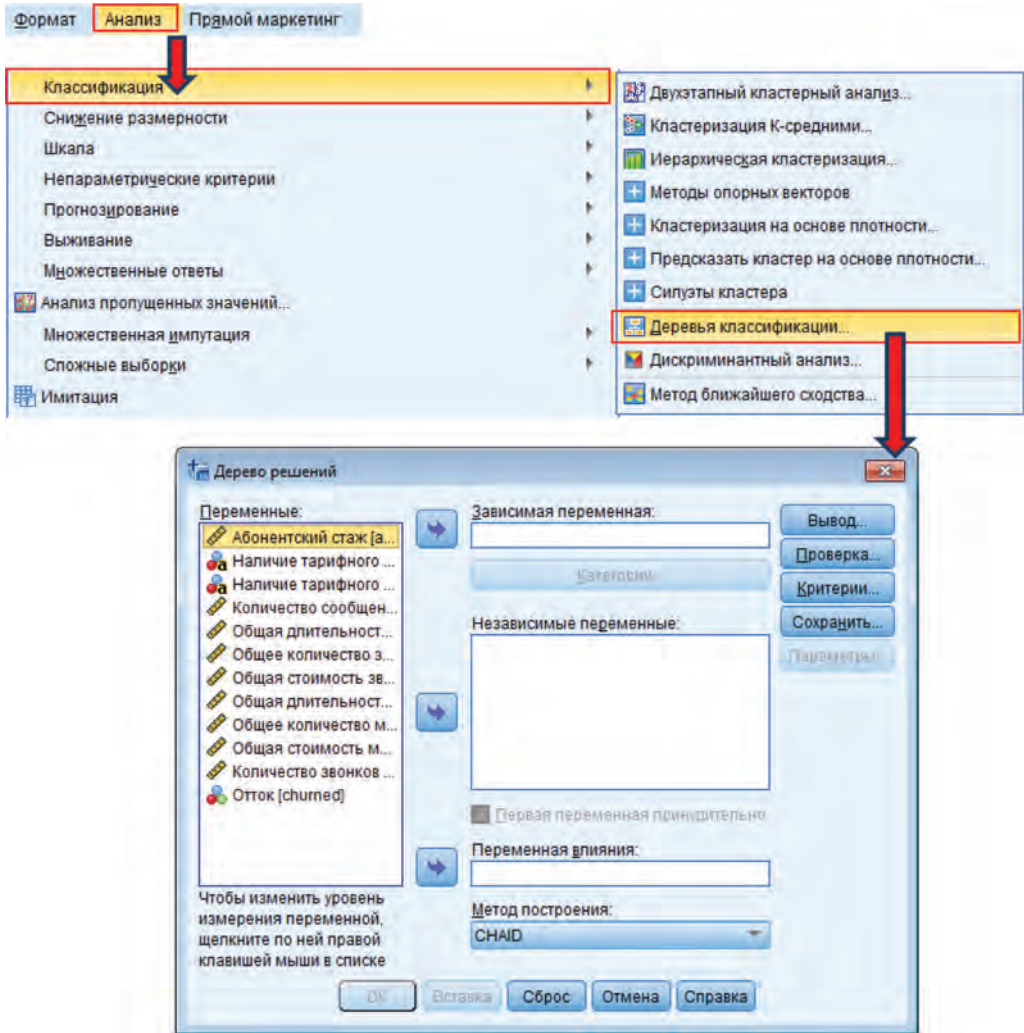


Рис. 2.1 ❖ Запуск процедуры **Дерева классификации**

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**, **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.

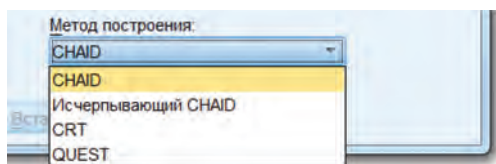


Рис. 2.2 ❖ Выпадающий список **Метод построения**

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 шкале, 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 Биггом, де Виллем и 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 шкале. Предикторы могут быть измерены в номинальной, порядковой и количественной шкале (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 вырастить 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 (эффект переобучения). Для упрощения структуры дерева и устранения переобучения в методе CRT предусмотрена возможность отсекаания ветвей (прунинг). Прунинг позволяет получить дерево «подходящего размера», избежать построения ветвистых, усложненных деревьев и при этом достичь наиболее точной оценки классификации.

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 большего времени 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 шкале.

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, номинальная, порядковая.

Количественная шкала содержит информацию о расстояниях между уровнями переменной, порядке уровней и количестве объектов в уровнях. Пример предиктора с количественной шкалой – переменная *Возраст*. Например, я знаю, что расстояние между 25 и 30 в два раза меньше, чем расстояние между 30 и 40, 30-летний на 5 лет старше 25-летнего. Я могу упорядочить уровни по нарастанию или убыванию интенсивности определенного признака (например, по увеличению возраста): после 25 сле-