

Двумерный анализ

Loginom Хакатон 2020. Секции 1, 3

Введение

В прошлом выпуске мы показали, как строить модель предрасположенности клиента совершить определенное действие (на примере повторного визита).

Но мы подробно не остановились на важной теме – **двумерном анализе**, а по сути, **оптимальном квантовании**, которое обязательно перед применением логистической регрессии и реализовано в компоненте **Конечные классы** платформы Loginom.

Двумерный анализ

В предсказательной аналитике уделяют большое внимание анализу связей между переменными. Самым простым и типичным является случай анализа взаимосвязи (сопряженности) двух переменных, так называемый **ДВУМЕРНЫМ АНАЛИЗОМ**.

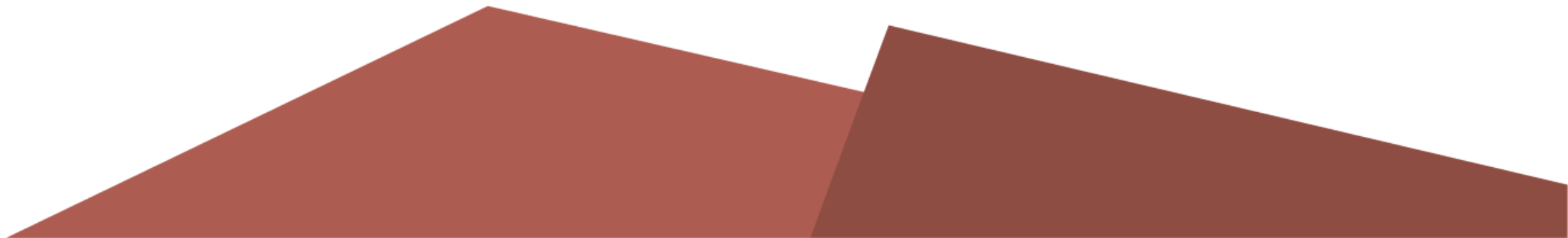
Формирование конечных классов (англ: Fine&Coarse Classing) как процесс оптимального квантования или биннинга (для непрерывных переменных) и снижения разнообразия категорий путем их объединения (для категориальных переменных) позволяет построить устойчивые модели оценки вероятности наступления события даже при использовании методов, линейных по своей сути (например, **логистическая регрессия**).

Двумерный анализ

Кроме того, процедуры формирования конечных классов позволяют получить два важных результата:

- оценить предсказательную силу отдельной переменной (анализ IV);
- оценить характер связи значения переменной с бинарной выходной переменной (анализ WoE).

Математика WoE-анализа



Первый шаг

Первый шаг в анализе взаимоотношений двух переменных является их перекрестная классификация, или построение таблицы сопряженности. Пример: реакция клиентов банка на предложение открыть кредитную карту (см. исходные данные в файле **Отклики.xls**).

Отклик	Уровень образования		
		низкий	высокий
	Нет (0)	50	310
	Да (1)	95	45
	Всего	145	355

Второй шаг

Методика анализа двумерной таблицы:

- Пересчитать абсолютные частоты в проценты;
- Сравнить процентные показатели, полученные для подгрупп с разным уровнем независимой переменной, каждый раз внутри одной категории зависимой переменной.

Отклик	Уровень образования			
		низкий	высокий	Всего
	Нет (0)	34,5%	87,3%	72%
	Да (1)	65,5%	12,7%	28%
	Всего	100%	100%	100%

WoE-анализ

WoE-анализ (Weight Of Evidence) или **совокупность доказательств** - статистический метод оценки влияния тех или иных факторов на справедливость некоторой гипотезы.

$$WoE_i = \ln \frac{N_i/N}{P_i/P}$$

где i – индекс признака, для которого вычисляется показатель WoE, N_i – число не-событий в группе, N – общее число не-событий, P_i – число событий в группе, P – общее число событий.

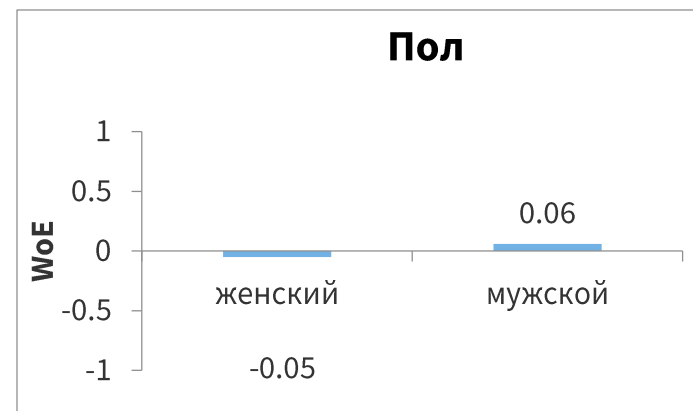
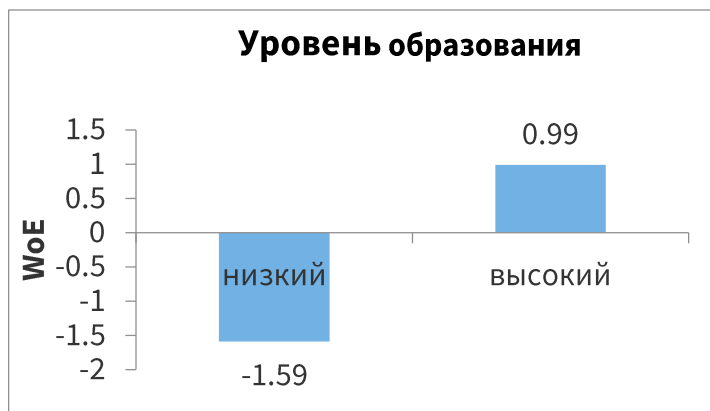
Гипотеза независимого поведения признаков: пропорция события и не-события в анализируемой подгруппе должна сохраняться такой же, как и для всей выборки в целом.

Интерпретация WoE-индексов

Если значение категории совпадает с событием большее число раз, чем с не-событием, то согласно формуле, под знаком логарифма будет значение меньше 1, что делает его отрицательным.

WoE < 0 указывает на большую вероятность появления события, а **WoE > 0** - не-события.

Индекс WoE есть **количественная мера предсказательной силы отдельной категории внутри переменной.**



Информационный индекс

WoE является промежуточным элементом для вычисления агрегированной величины, называемой информационным **индексом IV** (Information Value):

$$IV = \sum_{i=1}^K \left\{ \left(\frac{N_i}{N} - \frac{P_i}{P} \right) \cdot WoE_i \right\} ,$$

Информационный индекс отвечает за **предсказательную способность** (силу) всей переменной.

$$\begin{aligned} IV(\text{Уровень образования}) &= \left(\frac{50}{360} - \frac{95}{140} \right) \cdot (-1,59) + \left(\frac{310}{360} - \frac{45}{140} \right) \cdot (0,99) \\ &= 1,39. \end{aligned}$$

Значимость

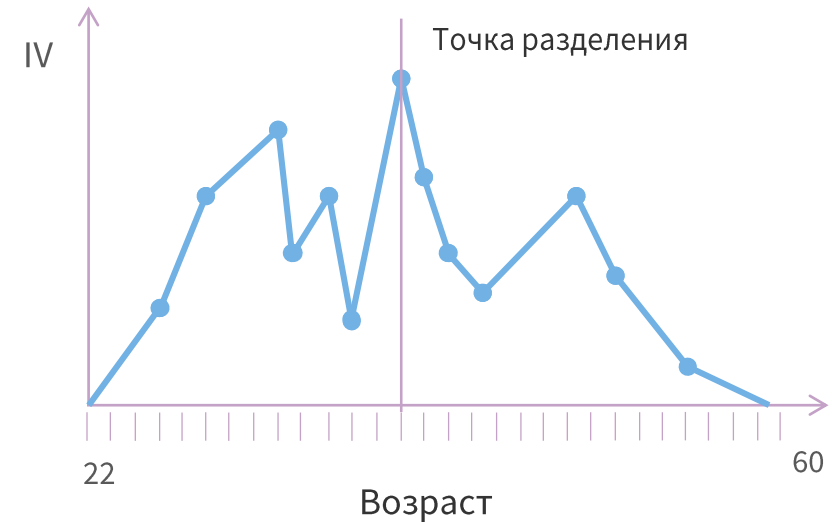
С помощью значения IV можно оценить значимость (предсказательную силу) переменной с бинарной выходной переменной:

- $IV < 0,02$ - отсутствует;
- $0,02 \leq IV < 0,1$ - низкая;
- $0,1 \leq IV < 0,3$ - средняя;
- $IV \geq 0,3$ - высокая.

На практике предсказание отклика по одной, двум значимым переменным, сомнительна. При решении реальных задач приходится иметь дело с большим числом категориальных переменных, имеющих большее число уникальных значений, а также с непрерывными переменными.

Обработка непрерывной переменной

Идеи конечных классов для обработки непрерывной переменной те же самые. На первом шаге все уникальные значения переменной сортируются по возрастанию, и для каждого из них рассчитывается WoE-значение, а для всей переменной – индекс IV (начальные классы). На следующем шаге производится итерационное формирование конечных классов путем присоединения соседних начальных классов последовательно начиная с первого с новым расчетом IV.



По сути, на каждой итерации выполняется поиск оптимальной точки разделения, на основе которой будут созданы два конечных класса. Это точка, в которой значение IV будет максимальным.

Процесс выполняется итерационно до тех пор, пока не будут достигнуты ограничения на максимально допустимое число конечных классов и минимальную долю примеров в классе.

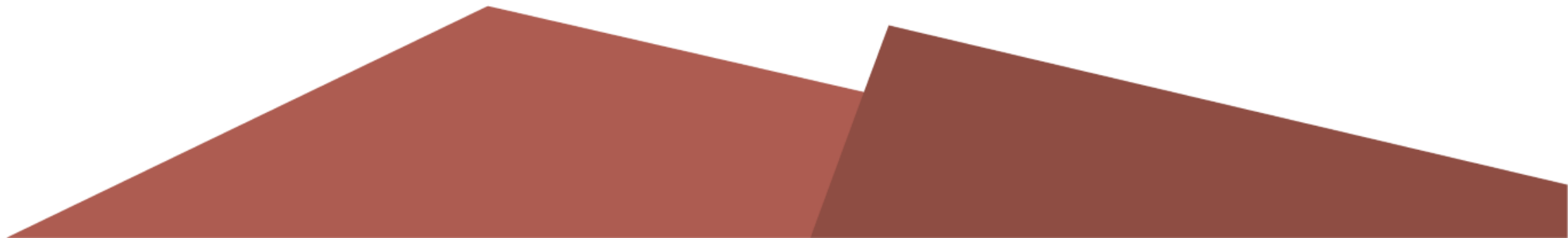
Применение

- Эффективное использование логистической регрессии.
- Переменные с отсутствующей значимостью можно не подавать на вход логистической регрессии.
- Из нескольких коррелирующих между собой переменных можно оставить одну с максимальным значением IV.

Пример

Пример расчетов можно посмотреть в файле **woe-анализ.xls**.

Работа в Loginom



Файлы

Для дальнейшего понимания откройте сценарий **woe-анализ.lgp**.

Сценарий

Пакеты woe_анализ Модуль1 Сценарий

Компоненты

лог

Data Mining

Логистическая регрессия

Отклики

Конечные классы

Логистическая регрессия

Производные компоненты +

Подкл

Откройте сценарий и визуализатор Куб первого узла.

Куб

Пакеты

woe_анализ

Модуль1

Сценарий

Отклики

Визуализаторы

Куб

+

Уровень образования

+

Σ Факты

Отклик

+

	низкий	высокий	Итого:
Ложь	50	310	360
Истина	95	45	140
Итого:	145	355	500

Наблюдаем таблицу перекрестной классификации.

Куб

Пакеты

woe_анализ

Модуль1

Сценарий

Отклики

Визуализаторы

Куб

+

Уровень образования

+

Σ Факты

Отклик

+

	низкий	высокий	Итого:
Ложь	50	310	360
Истина	95	45	140
Итого:	145	355	500

Наблюдаем таблицу перекрестной классификации.

Настройка

Пакеты

woe_анализ

Модуль1

Сценарий

Конечные классы

Настройка

Настройка назначений столбцов

Состояние входа

Не активировано

Активировать

№	Входные поля	Вид данных	<input checked="" type="checkbox"/>	Назначение
1	ab Уровень образования	Дискретный	<input checked="" type="checkbox"/>	Входное
2	ab Пол	Дискретный	<input checked="" type="checkbox"/>	Входное
3	0/1 Отклик	Дискретный	<input checked="" type="checkbox"/>	Выходное

Назначение столбца

Выходное

Настройки выходного поля

Настраиваемое значение "события"true

Назад

Далее

Далее берем узел Конечные классы и настраиваем столбцы: Уровень образования, Пол – входные, Отклик – выходное, событие - true

Настройка

Пакетыwoe_анализМодуль1СценарийКонечные классыНастройка

Настройка назначений столбцов

Состояние входаНе активировано

Активировать

№	Входные поля	Вид данных	<input checked="" type="checkbox"/>	Назначение
1	ab Уровень образования	Дискретный	<input checked="" type="checkbox"/>	Входное
2	ab Пол	Дискретный	<input checked="" type="checkbox"/>	Входное
3	0/1 Отклик	Дискретный	<input checked="" type="checkbox"/>	Выходное

Назначение столбца

Выходное

Настройки выходного поля

Настраиваемое значение "события"true

Назад

Далее

Далее берем узел Конечные классы и настраиваем столбцы: Уровень образования, Пол – входные, Отклик – выходное, событие - true

Настройка

Пакетыwoe_анализМодуль1СценарийКонечные классыНастройка

Настройка конечных классов

Состояние входаВход активированАктивировано

ab Уровень образования1,39

ab Пол0,00

Конечные классы

Минимальный вес, %5

Максимальное кол-во5

Установить...

Оптимизация

Равномерность, %0

Уровень образования

низкийвысокий

ДоляWoEIV

№	Метка	Нижняя	Верхняя	События	Не-события	Всего	Доля	Вес доказательства	Инф.индекс
0	[низкий]		низкий, высо...	95	50	145	<div><div></div>29%</div>	-1,59	0,86
1	[высокий]	низкий, высо...		45	310	355	<div><div></div>71%</div>	0,99	0,53

1.38793596177...

На следующем шаге мы видим диаграмму WoE-индексов, таблицу с показателями и значения IV для переменных. Здесь можно вмешаться в автоматический расчет и сдвинуть границы в любую сторону.