

Данные из нескольких таблиц

Данные из нескольких таблиц

Семейный статус ДОЛЖНИКОВ

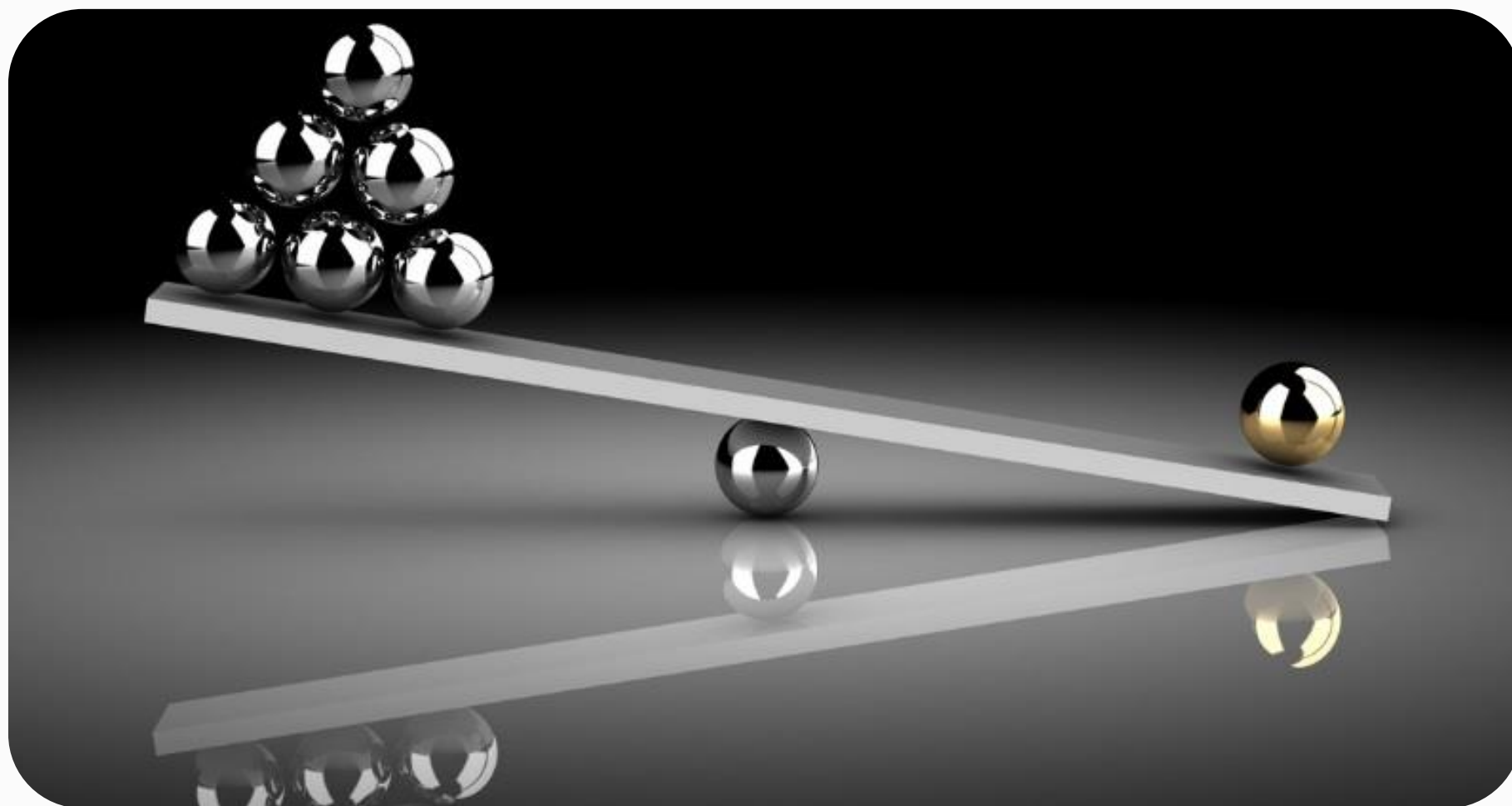
married	tertiary	1101	yes	no	yes
married	secondary	22	no	no	yes
married	secondary	1466	yes	no	yes
married	tertiary	473	no	no	yes
married	primary	0	no	no	yes
married	secondary	521	yes	no	yes
married	primary	2103	yes	no	yes
married	secondary	986	no	yes	yes
married	secondary	-157	no	yes	yes
married	primary	512	yes	yes	yes
married	tertiary	1612	no	no	yes
married	secondary	-233	yes	no	yes
married	secondary	71	no	no	yes
married	unknown	1	no	no	yes
married	secondary	387	yes	no	yes
married	tertiary	0	yes	yes	yes
married	secondary	380	no	no	yes
married	secondary	-757	yes	yes	yes
married	secondary	23495	no	no	yes
married	tertiary	3634	no	no	yes
married	secondary	6	no	no	yes
married	secondary	2	yes	no	yes
married	primary	0	no	no	yes
married	primary	3777	yes	no	yes
married	secondary	614	no	no	yes
married	tertiary	28	yes	no	yes

Данные из нескольких таблиц

Признаки используются
для предсказания
целевой переменной.

Данные из нескольких таблиц

Дисбаланс классов



Данные из нескольких таблиц

Что же делать?

- 1 Заменить на значение, которое встречается чаще остальных

Данные из нескольких таблиц

Что же делать?

- 1 Заменить на значение, которое встречается чаще остальных
- 2 Удалить строки с пропусками

Данные из нескольких таблиц

Что же делать?

- 1 Заменить на значение, которое встречается чаще остальных
- 2 Удалить строки с пропусками
- 3 Найти дополнительный источник данных

Данные из нескольких таблиц

VLOOKUP

Индекс 1

Индекс 2

A2 – начало диапазона

Направление поиска

Критерий: ищем число 39

	Индекс 1	Индекс 2	
1	id	default	
2	7	no	
3	12	no	
4	13	no	
5	14	no	
6	24	no	<div>=?=VLOOKUP(39; A2:B15; 2; FALSE)</div>
7	25	yes	
8	27	no	
9	29	no	
10	31	no	
11	39	yes	
12	44	no	
13	46	no	
14	47	no	
15	49	yes	
16	50	yes	

Нужно значение этой ячейки

B15 – конец диапазона