

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

Задача кредитного скоринга

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

Ошибка

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

Некорректное значение

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23	?	40000	А	1
Ж	30	Екб	19973	В	1

Пропущенное значение

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

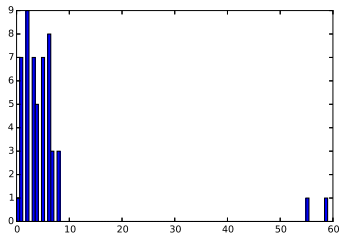
Выброс? Ошибка? Специальное значение?

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

Подозрительная точность

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	100000000	С	1
М	10	Мск	40000	С	0
А	28	СПб	80000	В	0
М	23		40000	А	1
Ж	30	Екб	19973	В	1

Что означают эти буквы? Можно ли их сравнивать?



Статистики:

- Первая квартиль  $Q_1$ :  $\frac{1}{\ell} \sum_{i=1}^{\ell} [x_i \leq Q_1] = 0.25$
- Третья квартиль  $Q_3$ :  $\frac{1}{\ell} \sum_{i=1}^{\ell} [x_i \leq Q_3] = 0.75$
- Интерквартильный размах:  $IQR = Q_3 - Q_1$

Эвристика: выбросы лежат за пределами отрезка  $[Q_1 - 1.5 IQR, Q_3 + 1.5 IQR]$ .



Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	?	С	1
М	10	?	40000	С	0
?	28	СПб	80000	В	0
М	23	?	40000	А	1
Ж	30	Екб	19973	В	1

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	49994.6	С	1
М	10	?	40000	С	0
?	28	СПб	80000	В	0
М	23	?	40000	А	1
Ж	30	Екб	19973	В	1

Пропуски в числовых признаках:

- Замена на среднее

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	40000	С	1
М	10	?	40000	С	0
?	28	СПб	80000	В	0
М	23	?	40000	А	1
Ж	30	Екб	19973	В	1

Пропуски в числовых признаках:

- Замена на среднее
- Замена на медиану

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	?	С	1
М	10	Мск	40000	С	0
М	28	СПб	80000	В	0
М	23	Мск	40000	А	1
Ж	30	Екб	19973	В	1

Пропуски в категориальных признаках:

- Замена на самое популярное значение

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	?	С	1
М	10	Пр	40000	С	0
Пр	28	СПб	80000	В	0
М	23	Пр	40000	А	1
Ж	30	Екб	19973	В	1

Пропуски в категориальных признаках:

- Замена на самое популярное значение
- Замена на новое значение

Пол	Возраст	Город	Доход	Образование	Вернул?
М	25	Мск	70000	А	1
Ж	31	Мск	?	С	1
М	10	Екб	40000	С	0
Ж	28	СПб	80000	В	0
М	23	Мск	40000	А	1
Ж	30	Екб	19973	В	1

Пропуски в категориальных признаках:

- Замена на самое популярное значение
- Замена на новое значение
- Случайный выбор из распределения значений

- В данных может быть много проблем: пропуски, выбросы, некорректные значения
- Выбросы можно обнаруживать с помощью интерквартильного размаха
- Пропуски можно исправлять путем замены на разумные значения