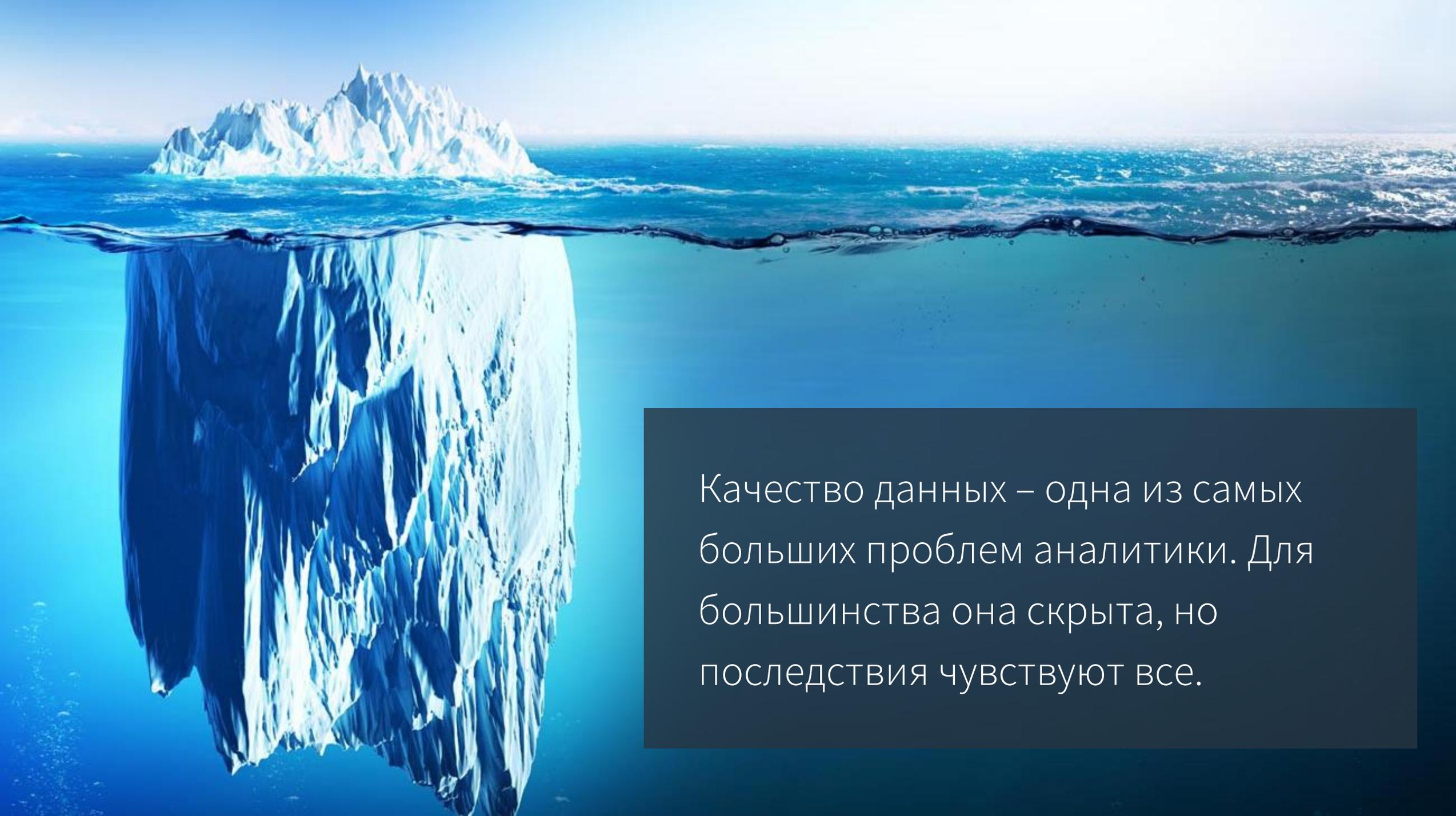


Визуализация качества данных в Logiном

Арустамов Алексей



Качество данных – одна из самых больших проблем аналитики. Для большинства она скрыта, но последствия чувствуют все.

Критерии качества

1. Согласованность — взаимная непротиворечивость
2. Своевременность — доступность в нужный момент
3. Целостность — корректность ссылок, соответствие правилам
4. Точность — нужный уровень детализации
5. Полнота — достаточность объема, глубины и широты
6. Правильность — отсутствие систематических ошибок
7. ...

Способы повышения качества

1. **Профилирование** – исследование данных для выяснения их статистических характеристик, таких как распределения величин, наличие выбросов, параметры выборки
2. **Верификация** – проверка данных по критериям достоверности источника, точности, согласованности и соответствия формату
3. **Стандартизация** – приведение к определённому формату и представлению, который обеспечивает корректный анализ

Визуализатор «Качество данных»

Расчет характеристик выборок для выявления типичных индикаторов проблем в данных:

1. Пропуски и пробельные значения
2. Выбросы и экстремальные значения
3. Наличие моды и монотонность
4. Распределения значений
5.

№	Тип	Метка	Вид	Проблемы	Виды проблем	Характеристика набора данных	Значение
6	ab	Группа клиента	⚙	0.02%	Пустые - 0,02% (15)	Метод определения нетипичных значений	Стандартное отклонение
13	90	Цена за единицу	⚙	0.01%	Пропуски - 0,00% (2) Нули - 0,01% (8)	Столбцов	16
15	90	Сумма с учетом скидки	⚙	0.01%	Нули - 0,01% (10)	Строк	99 736
9	ab	Наименование товара	⚙	0.01%	Пустые - 0,01% (5)	Заполненных полей	87,50%
0	12	ID	⚙	0.00%	Пропуски - 0,00% (4)	Полных записей	99,99%
3	ab	Город	⚙	0.00%	Пробелы - 0,00% (3)	Пригодных столбцов	15 из 16
7	ab	Отдел	⚙	0.00%	Константа	Индекс EPV	6 233,50
1	11	Дата продажи	⚙	✓			
2	ab	№ Клиента	⚙	✓			
4	ab	Экономический район	⚙	✓			
5	ab	Федеральный округ	⚙	✓			
8	12	Артикул	⚙	✓			
10	ab	Группа товара	⚙	✓			
11	ab	Обобщенная группа товаров	⚙	✓			
12	ab	Единица измерения	⚙	✓			
14	12	Количество	⚙	✓			

Сводная оценка полей с
указанием выявленных проблем

Качество данных

Сводка | Дискретные | Непрерывные | Показатели | XLS

№	Тип	Метка	Индекс качества	Пропуски	Гистограмма	Значения	Экстремальные	Выбросы	Пустые	Пробельные	Пробелы в конце	Длины строк	Диапазон значений	Мода	Монотонность	
2	ab	№ Клиента	72%	0			455	1 871	0	0	0	1 — 8	10...EX000005	9999...	Не монотонно	
3	ab	Город	67%	0			0	2 047	0	3	0	2 — 24	Пробелы...Ярославль	Мос...	Не монотонно	
4	ab	Экономический район	67%	0			0	522	0	0	0	17 — 37	Волго-Вятский экономический район...	Цен...	Не монотонно	
5	ab	Федеральный округ	64%	0			0	2 395	0	0	0	23 — 34	Дальневосточный федеральный округ...	Цен...	Не монотонно	
6	ab	Группа клиента	67%	0			15	0	15	0	0	0 — 17	Пустое...Постоянный клиент	Кли...	Не монотонно	
7	ab	Отдел	0%	0			0	0	0	0	0	21 — 21	Корпоративные клиенты...Корпорати...	Кор...	Не монотонно	
8	ab	Артикул	93%	0			0	1 446	0	0	0	6 — 6	100005...113822	102439	Не монотонно	
9	ab	Наименование товара	93%	0			0	1 454	5	0	0	3	0 — 105	Пустое...Эмаль НЦ-132П OLEKOLOR ч...	Пли...	Не монотонно
10	ab	Группа товара	80%	0			2	384	0	0	0	2	5 — 39	Армирующие материалы...Эмаль	Стен...	Не монотонно
11	ab	Обобщенная группа то	73%	0			102	0	0	0	0	10 — 44	Армирующие материалы...Сыпучие и	Отд...	Не монотонно	
12	ab	Единица измерения	70%	0			102	0	0	0	0	1 — 4	кв.м...шт	шт	Не монотонно	

Формат | Сортировка | Найти | XLS

#	Федеральный округ	ab Группа клиента	ab Отдел	ab Артикул	ab Наименование товара	ab Группа товара	ab Обобщенная группа товаров	ab Единица измерения	9.0 Цена за единицу	12 Кол
1	И федеральный округ	Постоянный клиент	Корпоративные клиенты	102341	Профиль угловой арочный стальной О, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
2	И федеральный округ	Постоянный клиент	Корпоративные клиенты	102342	Профиль угловой арочный стальной Р, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
3	И федеральный округ	VIP клиент	Корпоративные клиенты	110485	Пароизоляция ИЗОСПАН В 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	9,58	
4	И федеральный округ	Клиент	Корпоративные клиенты	102341	Профиль угловой арочный стальной О, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
5	И федеральный округ	VIP клиент	Корпоративные клиенты	110488	Гидро-пароизоляция ИЗОСПАН С 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	14,6	
6	И федеральный округ	VIP клиент	Корпоративные клиенты	110491	Гидро-пароизоляция ИЗОСПАН D универсальная 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	13,21	
7	И федеральный округ	VIP клиент	Корпоративные клиенты	110492	Гидро-ветрозащитная ИЗОСПАН АМ 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	20,33	
8	И федеральный округ	VIP клиент	Корпоративные клиенты	110496	Гидро-пароизоляция ИЗОСПАН FD отражающая 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	22,33	
9	И федеральный округ	Клиент	Корпоративные клиенты	110488	Гидро-пароизоляция ИЗОСПАН С 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	14,6	
10	И федеральный округ	Клиент	Корпоративные клиенты	102340	Профиль угловой арочный стальной N, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
11	И федеральный округ	Клиент	Корпоративные клиенты	102342	Профиль угловой арочный стальной Р, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
12	И федеральный округ	VIP клиент	Корпоративные клиенты	102341	Профиль угловой арочный стальной О, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
13	И федеральный округ	VIP клиент	Корпоративные клиенты	102342	Профиль угловой арочный стальной Р, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
14	И федеральный округ	VIP клиент	Корпоративные клиенты	110488	Гидро-пароизоляция ИЗОСПАН С 70 кв.м., кв.м.	Гидроизоляция	Изоляционные материалы	кв.м	14,6	
15	И федеральный округ	Клиент	Корпоративные клиенты	102340	Профиль угловой арочный стальной N, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	
16	И федеральный округ	VIP клиент	Корпоративные клиенты	102341	Профиль угловой арочный стальной О, 2,75 м, м	Профиль заказной	Металлический профиль и комплектующие	м	6,47	

384

Потенциальные проблемы
в дискретных полях

Главная Вставка Разметка страницы Формулы Данные Рецензирование Вид Настройки Справка Что вы хотите сделать?

Вставить Вырезать Копировать Формат по образцу Буфер обмена

Шрифт Calibri 10 A A Ж К Ц

Выравнивание Переносить текст Объединить и поместить в центре

Число Текстовый % 000 ,00 ,00

Условное форматирование Форматировать как таблицу

Стили Обычный Нейтральный Плохой Хороший

Ячейки Вставить Удалить Формат

Редактирование Автосумма Заполнить Очистить Сортировка и фильтр Найти и выделить

C2 ID-клиента

	C	D	E	F	G	H	I	J	K	L	M	N
1	Метка	Индекс качества	Пропуски	Экстремальные	Выбросы	Пустые	Пробельные	Пробелы в конце	Длины строк	Диапазон значений	Мода	Монотонность
2	ID-клиента	1,00	0	0	0	Недоступно	Недоступно	Недоступно	Недоступно	2 476 ... 83 467	2476	Не монотонная
3	Гараж	0,99	0	0	0	54	0	0	0 — 3	... Нет	Да	Не монотонная
4	Давать кредит(число)	0,97	0	0	0	Недоступно	Недоступно	Недоступно	Недоступно	0 ... 1	0	Не монотонная
5	Давать кредит	0,97	0	0	0	0	0	0	2 — 3	Да ... Нет	Нет	Не монотонная
6	Срок эксплуатации машины	0,93	92	0	0	0	0	3	3 1 — 3	... INF	3	Не монотонная
7	Основное направление расходов	0,90	0	0	0	2	0	0	0 4 — 37	Выплаты по кредитам/займам ... Турпоездки, р	Одежда, продукты питания и т.п.	Не монотонная
8	Цель кредитования	0,89	0	0	0	0	0	0	0 4 — 33	Иное ... Турпоездки, развлечения и т.п.	Покупка товара	Не монотонная
9	Цель кредитования	0,86	0	0	0	0	0	5	5 1 — 33	... Турпоездки, развлечения и т.п.	Покупка товара	Не монотонная
10	Гражданское состояние	0,86	0	0	0	0	0	0	0 2 — 3	Да ... Нет	Да	Не монотонная
11	Расположение	0,85	0	0	0	0	0	0	0 5 — 7	область ... центр	центр	Не монотонная
12	Загородный дом	0,85	0	0	0	0	0	90	90 1 — 3	... Нет		Не монотонная
13	Должность	0,82	0	0	0	0	0	0	0 11 — 13	неруководящая ... руководящая	неруководящая	Не монотонная
14	Машина	0,82	34	0	0	0	0	0	0 9 — 14	импортная ... отечественная	отечественная	Не монотонная
15	Образование	0,82	0	0	0	0	14	15	15 0 — 11	... специальное	специальное	Не монотонная
16	Способ приобретения собств.	0,79	0	0	0	0	0	0	0 6 — 25	другое ... покупка	другое	Не монотонная
17	Пол	0,74	0	0	5	0	0	5	5 1 — 3	... Муж	Муж	Не монотонная
18	Примечание	0,00	0	0	0	0	149	0	0 0 — 0	...		Не монотонная
19	ID Банка	0,00	0	0	0	0	0	0	0 6 — 6	B70989 ... B70989	B70989	Не монотонная

Выгрузка отчета в Excel

Визуализатор «Качество данных»

позволяет:

1. Быстро оценить пригодность выборки для анализа
2. Автоматически провести типовые проверки качества данных
3. Локализовать проблемы

loginom.ru