

Технология ETL

Вариант 2

а) определить отсутствующие данные (тепловая карта, процентный список, гистограмма с отсутствующими данными) и принять обоснованное решение о дальнейшей работе с ними (отбросить наблюдения, отбросить параметр, заменить отсутствующие данные)

```
# Тепловая карта пропущенных значений (первые 30 колонок)
cols = df.columns[:30]
colours = ['#000099', '#ffff00'] # синий = данные, желтый = пропуски
sns.heatmap(df[cols].isnull(), cmap=sns.color_palette(colours))
plt.title('Тепловая карта пропущенных значений (первые 30 колонок)')
plt.tight_layout()
plt.show()

# Процент пропусков по каждому признаку
print("Процент пропусков по признакам:")
dict_missing_percentages = {}
for col in df.columns:
    pct_missing = np.mean(df[col].isnull())
    if pct_missing > 0: # Показываем только колонки с пропусками
        print(f'{col} - {round(pct_missing*100)}%')
        dict_missing_percentages[col] = pct_missing

# Создаем индикаторы пропусков
for col in df.columns:
    if df[col].isnull().any():
        df[f'{col}_ismissing'] = df[col].isnull()

# Гистограмма количества пропусков на строку
ismissing_cols = [col for col in df.columns if 'ismissing' in col]
df['num_missing'] = df[ismissing_cols].sum(axis=1)

plt.figure(figsize=(10, 6))
df['num_missing'].value_counts().sort_index().plot.bar()
plt.title('Количество пропусков на строку')
plt.xlabel('Число пропусков')
plt.ylabel('Количество строк')
plt.tight_layout()
plt.show()

# Гистограмма количества пропусков по признакам
missing_counts = df[ismissing_cols].sum(axis=0)
missing_counts.index = missing_counts.index.str.replace('_ismissing', '')
plt.figure(figsize=(12, 6))
missing_counts.sort_values().plot.bar()
plt.ylabel('Количество пропусков')
```

```

plt.xlabel('Признаки')
plt.title('Количество пропусков по признакам')
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()

# Принятие решения о пропусках
print(f"Исходное количество строк: {len(df)}")

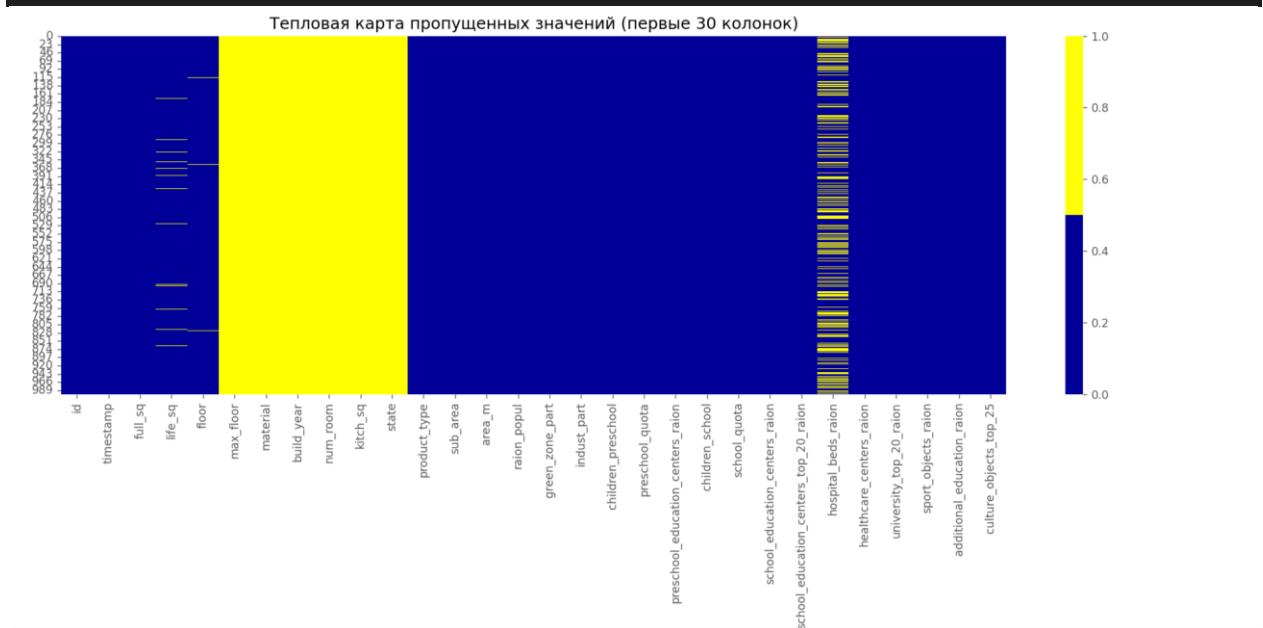
# Отбрасываем строки, где пропусков слишком много (>2)
ind_missing = df[df['num_missing'] > 2].index
df = df.drop(ind_missing, axis=0)
print(f"Осталось строк после удаления строк с >2 пропусками: {len(df)}")

# Удаляем признак с наибольшим числом пропусков (если процент пропусков > 50%)
if dict_missing_percentages:
    col_with_max_missing = max(dict_missing_percentages,
key=dict_missing_percentages.get)
    if dict_missing_percentages[col_with_max_missing] > 0.5:
        df = df.drop(columns=[col_with_max_missing])
        print(f"Удален столбец с наибольшим количеством пропусков:
{col_with_max_missing}")

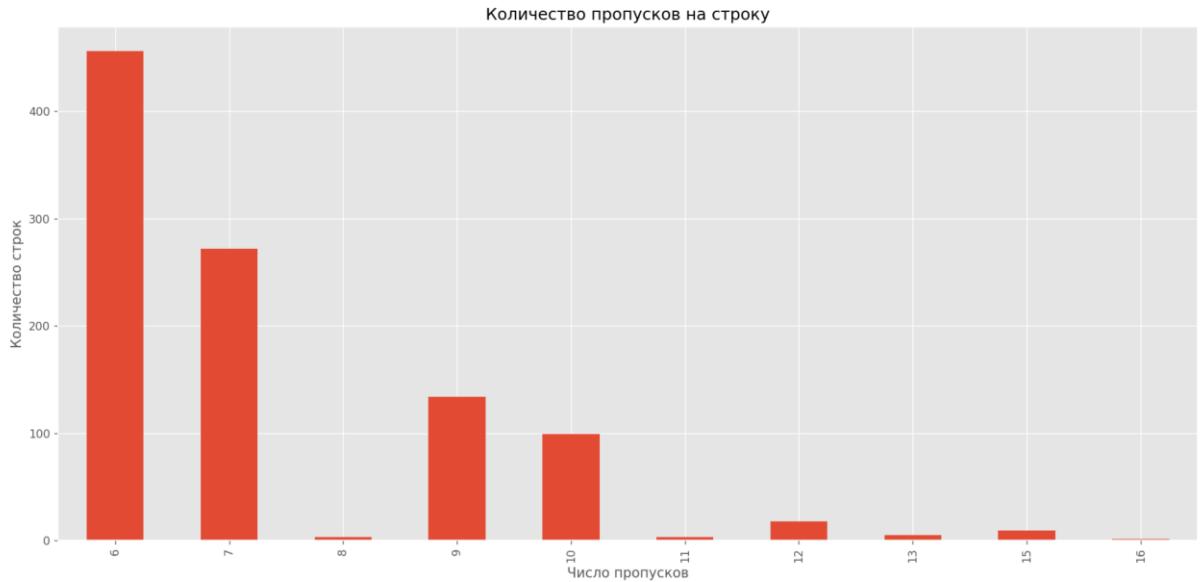
# Замена пропусков (числовые - медианой, категориальные - модой)
for col in df.select_dtypes(include=[np.number]):
    if df[col].isnull().any():
        df[col] = df[col].fillna(df[col].median())

for col in df.select_dtypes(exclude=[np.number]):
    if df[col].isnull().any():
        df[col] = df[col].fillna(df[col].mode()[0])

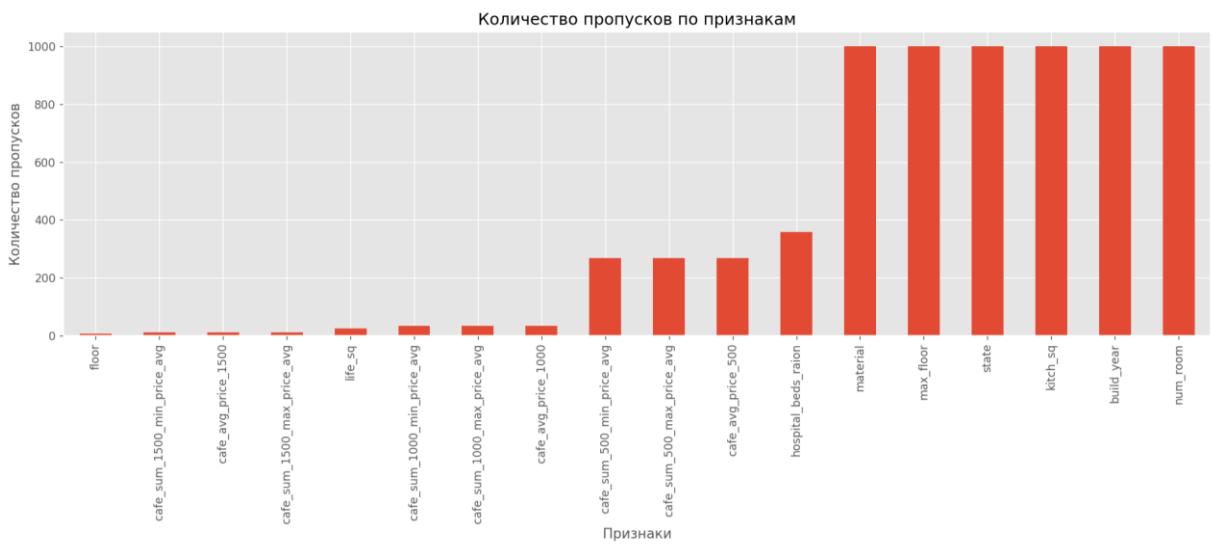
```



Тепловая карта показывает распределение пропусков. Она двухцветная: синий(данные присутствуют) и желтый (отсутствуют). По оси x – названия признаков, по оси у – номера строк наблюдения. Есть столбцы, где для второй тысячи строк из таблицы данные отсутствуют.



По оси x – количество пропущенных значений в одной строке, по оси у – количество строк с соответствующим числом пропусков.



По оси x – названия признаков, по оси у – количество пропущенных значений в каждом признаке. Помогает понять, какие признаки не требуют вмешательства.

```
life_sq - 2%
floor - 1%
```

```
max_floor - 100%
material - 100%
build_year - 100%
num_room - 100%
kitch_sq - 100%
state - 100%
hospital_beds_raion - 36%
cafe_sum_500_min_price_avg - 27%
cafe_sum_500_max_price_avg - 27%
cafe_avg_price_500 - 27%
cafe_sum_1000_min_price_avg - 3%
cafe_sum_1000_max_price_avg - 3%
cafe_avg_price_1000 - 3%
cafe_sum_1500_min_price_avg - 1%
cafe_sum_1500_max_price_avg - 1%
cafe_avg_price_1500 - 1%
```

Осталось строк после удаления: 600

Удален столбец с наибольшим количеством пропусков: hospital_beds_raion

б) определить наличие выбросов в данных (построение гистограммы, использование описательной статистики) и принять обоснованное решение о дальнейшей работе с ними (отбрасываем, корректируем, оставляем)

```
numeric_cols = df.select_dtypes(include=[np.number]).columns

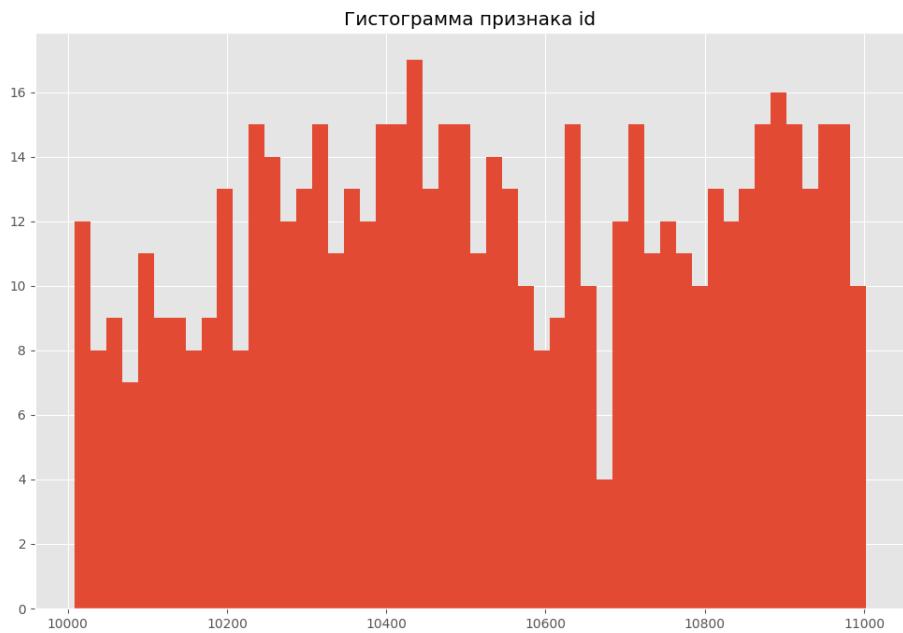
# Гистограммы для первых 5 числовых признаков
for col in numeric_cols[:5]:
    plt.figure(figsize=(10, 6))
    df[col].hist(bins=50)
    plt.title(f'Гистограмма признака {col}')
    plt.xlabel(col)
    plt.ylabel('Частота')
    plt.tight_layout()
    plt.show()

# Коробчатые диаграммы
for col in numeric_cols[:5]:
    plt.figure(figsize=(10, 6))
    df.boxplot(column=[col])
    plt.title(f'Boxplot признака {col}')
    plt.tight_layout()
    plt.show()

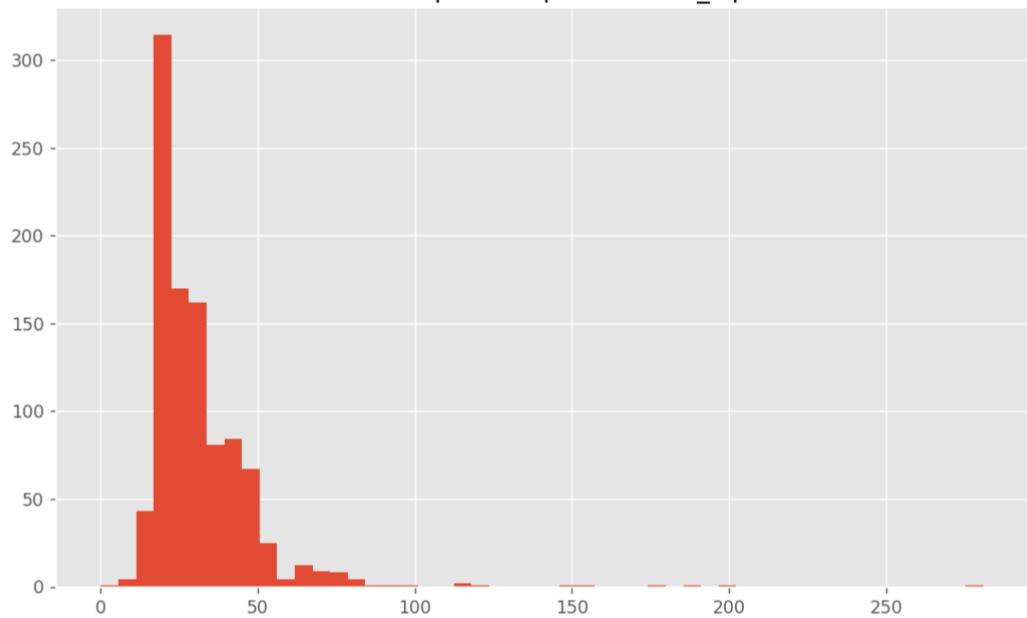
# Описательная статистика
print("Описательная статистика числовых признаков:")
print(df[numeric_cols].describe())
```

Гистограммы (анализы распределения) показывают форму распределения, плотность данных в различных диапазонах, наличие разрывов данных.

Приведем в пример первые 3 признака и их гистограммы:

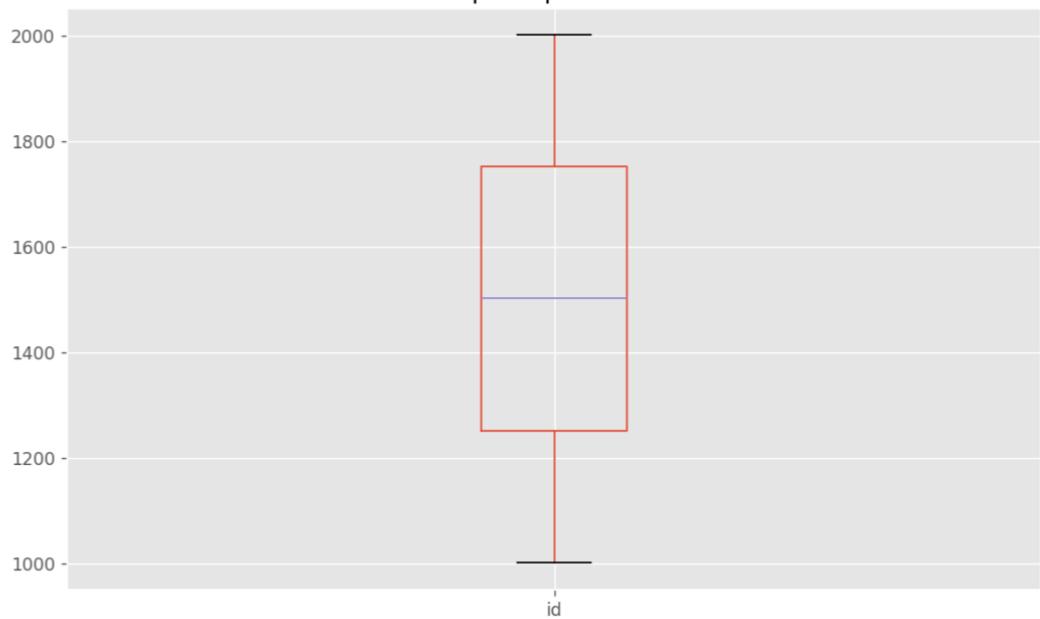


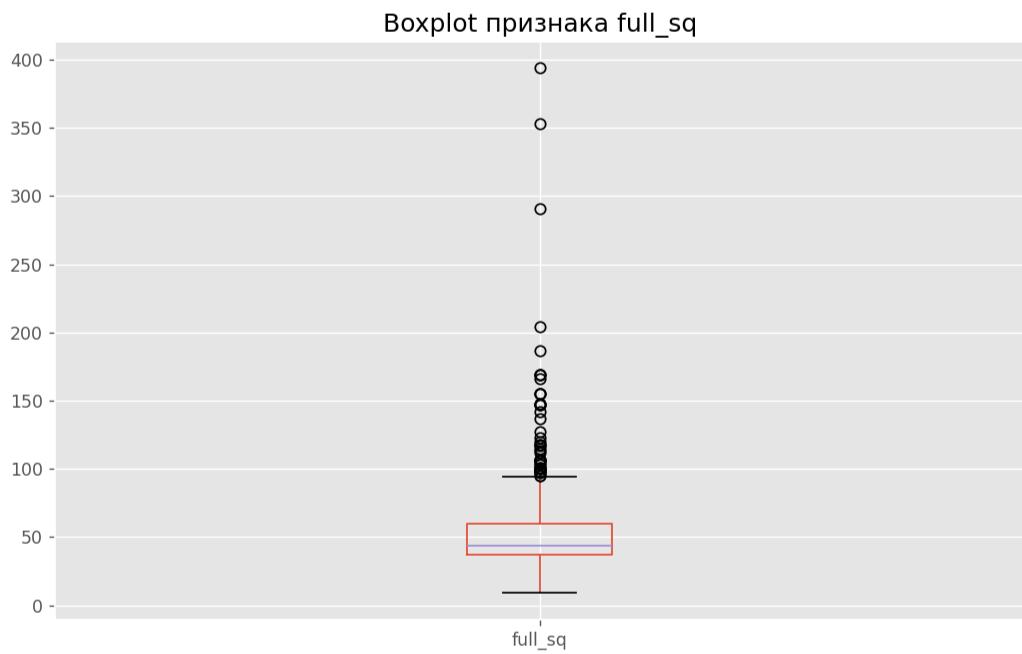
Гистограмма признака life_sq



Детекция выбросов показывает медиану и квартили распределения – 3 значения, которые делят упорядоченный набор на равные части. Выбросы – точки за пределами усов.

Boxplot признака id





в) определение ненужных данных (не добавляют ценности для решения данной задачи) - повторяющиеся данные, не имеют отношения к поставленной задаче, слишком много строк с одинаковыми значениями

```
# Поиск неинформативных признаков (более 95% одинаковых значений)
num_rows = len(df)
low_info_cols = []

for col in df.columns:
    cnts = df[col].value_counts(dropna=False, normalize=True)
```

```

if len(cnts) > 0:
    top_pct = cnts.iloc[0]
    if top_pct > 0.95:
        low_info_cols.append(col)
        print(f'{col}: {top_pct:.2%} одинаковых значений')

# Удаляем неинформативные столбцы
if low_info_cols:
    df = df.drop(columns=low_info_cols)
    print(f"Удалены неинформативные столбцы: {low_info_cols}")

# Удаляем дубликаты строк
initial_shape = df.shape
df = df.drop_duplicates()
final_shape = df.shape

if initial_shape[0] != final_shape[0]:
    print(f"Удалено дубликатов: {initial_shape[0] - final_shape[0]}")
print(f"После удаления дубликатов: {df.shape}")

```

```

max_floor: 100.00% одинаковых значений
material: 100.00% одинаковых значений
build_year: 100.00% одинаковых значений
num_room: 100.00% одинаковых значений
kitch_sq: 100.00% одинаковых значений
state: 100.00% одинаковых значений
product_type: 95.88% одинаковых значений
incineration_raion: 96.28% одинаковых значений
oil_chemistry_raion: 98.99% одинаковых значений
railroad_terminal_raion: 97.08% одинаковых значений
nuclear_reactor_raion: 95.77% одинаковых значений
water_1line: 95.37% одинаковых значений
big_road1_1line: 96.48% одинаковых значений
railroad_1line: 97.99% одинаковых значений
cafe_count_500_price_4000: 95.27% одинаковых значений
cafe_count_500_price_high: 98.49% одинаковых значений
mosque_count_500: 99.40% одинаковых значений
mosque_count_1000: 97.08% одинаковых значений
life_sq_ismissing: 97.99% одинаковых значений
floor_ismissing: 99.50% одинаковых значений

```

г) несогласованные - данные использование заглавных букв, форматы (время), адреса

```

# Нормализация строковых признаков (нижний регистр)
string_cols = df.select_dtypes(include=['object']).columns
for col in string_cols:
    df[col] = df[col].astype(str).str.lower().str.strip()

```

```
# Преобразование даты (если есть столбец timestamp)
if 'timestamp' in df.columns:
    df['timestamp'] = pd.to_datetime(df['timestamp'], errors='coerce')
    df['year'] = df['timestamp'].dt.year
    df['month'] = df['timestamp'].dt.month
    df['weekday'] = df['timestamp'].dt.weekday
```

Loginom

Текстовый файл:

Импорт из текстового файла

Имя файла / URL	<input type="text" value="birth-dates.txt"/>	...	+																																																																																																																														
Информация о файле	<input type="text" value="Нет"/>	Кодовая страница	<input type="text" value="Кириллическая (1251)"/>																																																																																																																														
Заголовок в первой строке	<input checked="" type="checkbox"/>	Пропустить строк	<input type="text" value="0"/>																																																																																																																														
<table border="1"> <tr><td>Дата рождения</td><td>Имя</td><td>Фамилия</td><td>Отчество</td><td>Код клиента</td></tr> <tr><td>05.10.1986</td><td>Август</td><td>Абайимов</td><td>Бенедиктович</td><td>1</td></tr> <tr><td>2-я-1990</td><td>Авдей</td><td>Абакумов</td><td>Богданович</td><td>2</td></tr> <tr><td>12.янв.90</td><td>Аверкий</td><td>Абакшин</td><td>Болеславович</td><td>3</td></tr> <tr><td>15 августа 1957</td><td>Аверьян</td><td>Абалакин</td><td>Бонифатович</td><td>4</td></tr> <tr><td>15 мая 1954</td><td>Аксентий</td><td>Абалаков</td><td>Бонифатиевич</td><td>5</td></tr> <tr><td>08.май.54</td><td>Автоном</td><td>Абалдуев</td><td>Борисович</td><td>6</td></tr> <tr><td>15.сен.56</td><td>Аган</td><td>Абалкин</td><td>Бориславович</td><td>7</td></tr> <tr><td>25 мак.59</td><td>Агафон</td><td>Абатурин</td><td>Брониславович</td><td>8</td></tr> <tr><td>15.12.1845</td><td>Аггей</td><td>Абатуров</td><td>Вавилич</td><td>9</td></tr> <tr><td>1947</td><td>Адам</td><td>Абашев</td><td>Вадимович</td><td>10</td></tr> <tr><td>1951 год</td><td>Андрян</td><td>Абашин</td><td>Валентинович</td><td>11</td></tr> <tr><td>1986 Feb 10</td><td>Питер</td><td>Азарий</td><td>Абашкин</td><td>Валерианович</td><td>12</td></tr> <tr><td>10.05.1986</td><td>Аким</td><td>Абаянцев</td><td>Валерьянович</td><td>13</td></tr> <tr><td>05.05.1986</td><td>Александр</td><td>Абдула</td><td>Валериевич</td><td>14</td></tr> <tr><td>19861005 Москва</td><td>Алексей</td><td>Абдулин</td><td>Варламович</td><td>15</td></tr> <tr><td>25.05.1986</td><td>Амвросий</td><td>Абдулов</td><td>Варламиевич</td><td>16</td></tr> <tr><td>04.05.1986</td><td>Амос</td><td>Абоимов</td><td>Варнавич</td><td>17</td></tr> <tr><td>5-10-1986 город Коломна</td><td>Ананий</td><td>Абраменко</td><td>Варсонофович</td><td>18</td></tr> <tr><td>05.05.1952</td><td>Анатолий</td><td>Абрамкин</td><td>Варсонофиевич</td><td>19</td></tr> <tr><td>05.01.1986</td><td>Андрей</td><td>Абрамов</td><td>Варфоломеевич</td><td>20</td></tr> <tr><td>05.10.1906</td><td>Андрон</td><td>Борищенко</td><td>Васильевич</td><td>21</td></tr> <tr><td>05.10.1986</td><td>Андроник</td><td>Борков</td><td>Вассианович</td><td>22</td></tr> <tr><td>05.10.1986</td><td>Аникей</td><td>Борковский</td><td>Велизарович</td><td>23</td></tr> <tr><td>дата 05101986</td><td>Аникита</td><td>Боровик</td><td>Велимирович</td><td>24</td></tr> </table>				Дата рождения	Имя	Фамилия	Отчество	Код клиента	05.10.1986	Август	Абайимов	Бенедиктович	1	2-я-1990	Авдей	Абакумов	Богданович	2	12.янв.90	Аверкий	Абакшин	Болеславович	3	15 августа 1957	Аверьян	Абалакин	Бонифатович	4	15 мая 1954	Аксентий	Абалаков	Бонифатиевич	5	08.май.54	Автоном	Абалдуев	Борисович	6	15.сен.56	Аган	Абалкин	Бориславович	7	25 мак.59	Агафон	Абатурин	Брониславович	8	15.12.1845	Аггей	Абатуров	Вавилич	9	1947	Адам	Абашев	Вадимович	10	1951 год	Андрян	Абашин	Валентинович	11	1986 Feb 10	Питер	Азарий	Абашкин	Валерианович	12	10.05.1986	Аким	Абаянцев	Валерьянович	13	05.05.1986	Александр	Абдула	Валериевич	14	19861005 Москва	Алексей	Абдулин	Варламович	15	25.05.1986	Амвросий	Абдулов	Варламиевич	16	04.05.1986	Амос	Абоимов	Варнавич	17	5-10-1986 город Коломна	Ананий	Абраменко	Варсонофович	18	05.05.1952	Анатолий	Абрамкин	Варсонофиевич	19	05.01.1986	Андрей	Абрамов	Варфоломеевич	20	05.10.1906	Андрон	Борищенко	Васильевич	21	05.10.1986	Андроник	Борков	Вассианович	22	05.10.1986	Аникей	Борковский	Велизарович	23	дата 05101986	Аникита	Боровик	Велимирович	24
Дата рождения	Имя	Фамилия	Отчество	Код клиента																																																																																																																													
05.10.1986	Август	Абайимов	Бенедиктович	1																																																																																																																													
2-я-1990	Авдей	Абакумов	Богданович	2																																																																																																																													
12.янв.90	Аверкий	Абакшин	Болеславович	3																																																																																																																													
15 августа 1957	Аверьян	Абалакин	Бонифатович	4																																																																																																																													
15 мая 1954	Аксентий	Абалаков	Бонифатиевич	5																																																																																																																													
08.май.54	Автоном	Абалдуев	Борисович	6																																																																																																																													
15.сен.56	Аган	Абалкин	Бориславович	7																																																																																																																													
25 мак.59	Агафон	Абатурин	Брониславович	8																																																																																																																													
15.12.1845	Аггей	Абатуров	Вавилич	9																																																																																																																													
1947	Адам	Абашев	Вадимович	10																																																																																																																													
1951 год	Андрян	Абашин	Валентинович	11																																																																																																																													
1986 Feb 10	Питер	Азарий	Абашкин	Валерианович	12																																																																																																																												
10.05.1986	Аким	Абаянцев	Валерьянович	13																																																																																																																													
05.05.1986	Александр	Абдула	Валериевич	14																																																																																																																													
19861005 Москва	Алексей	Абдулин	Варламович	15																																																																																																																													
25.05.1986	Амвросий	Абдулов	Варламиевич	16																																																																																																																													
04.05.1986	Амос	Абоимов	Варнавич	17																																																																																																																													
5-10-1986 город Коломна	Ананий	Абраменко	Варсонофович	18																																																																																																																													
05.05.1952	Анатолий	Абрамкин	Варсонофиевич	19																																																																																																																													
05.01.1986	Андрей	Абрамов	Варфоломеевич	20																																																																																																																													
05.10.1906	Андрон	Борищенко	Васильевич	21																																																																																																																													
05.10.1986	Андроник	Борков	Вассианович	22																																																																																																																													
05.10.1986	Аникей	Борковский	Велизарович	23																																																																																																																													
дата 05101986	Аникита	Боровик	Велимирович	24																																																																																																																													

Разбор даты рождения:

```
if ((day1 <> "", day1,
    RegExMatchedSubExp("\b(19|20)?[0-9]{2}[- /.,]([A-ZA-ЯЁ]{3,})[- /.,](0?[1-9]|12[0-9]|3[01])\b",
    Trim(Upper(birthdate)), 1)
)|
```

Выходные данные:

Разбор даты рождения — месяц задан строкой • Быстрый просмотр

#	ab d1	ab buthDay	ab m1	ab birthMonth	ab y1	ab birthYear	ab Дата рождения
1							05.10.1986
2							2-я-1990
3	12	12	ЯНВ	ЯНВ	90	90	12.янв.90
4	15	15	АВГУСТА	АВГУСТА	1957	1957	15 августа 1957
5	15	15	МАЯ	МАЯ	1954	1954	15 мая 1954
6	08	08					08.май.54
7	15	15	СЕН	СЕН	56	56	15.сен.56
8	25	25	МАК	МАК	59	59	25 мак 59
9							15.12.1845
10							1947
11							1951 год
12		10		FEB		19	1986 Feb 10 Питер
13							10.05.1986
14							05.05.1986
15							19861005 Москва
16							25.05.1986
17							04.05.1986
18							5-10-1986 город Коломна
19							05.05.1952
20							05.01.1986
21							05.10.1906
--							

Фильтр строк:

Фильтрация данных

Состояние входа

Не активировано

ab buthDay в списке	<input checked="" type="checkbox"/>	<input type="button" value="X"/>	<input type="button" value="+"/>
---------------------	-------------------------------------	----------------------------------	----------------------------------

Кодировка месяца:

Импорт из текстового файла

Имя файла / URL ... +

Информация о файле Кодовая страница

Заголовок в первой строке Пропустить строку

Месяц строковый	Первые 3 буквы	Номер месяца
ЯНВАРЬ	ЯНВ	01
ФЕВРАЛЬ	ФЕВ	02
МАРТ	МАР	03
АПРЕЛЬ	АПР	04
МАЙ	МАЙ	05
ИЮНЬ	ИЮН	06
ИЮЛЬ	ИЮЛ	07
АВГУСТ	АВГ	08
СЕНТЯБРЬ	СЕН	09
ОКТЯБРЬ	ОКТ	10
НОЯБРЬ	НОЯ	11
ДЕКАБРЬ	ДЕК	12
JANUARY	JAN	01
FEBRUARY	FEB	02
MARCH	MAR	03
APRIL	APR	04
MAY	MAY	05
JUNE	JUN	06
JULY	JUL	07
AUGUST	AUG	08
SEPTEMBER	SEP	09
OCTOBER	OCT	10
NOVEMBER	NOV	11
DECEMBER	DEC	12

Подсоединение из кодировки месяца

Настройка слияния данных

Тип операции

Фильтрация	
Столбцы основного набора данных	
<input checked="" type="checkbox"/>	ab birthMonth
<input checked="" type="checkbox"/>	ab birthYear
<input checked="" type="checkbox"/>	ab Дата рождения
<input checked="" type="checkbox"/>	ab Имя
<input checked="" type="checkbox"/>	ab Фамилия
<input checked="" type="checkbox"/>	ab Отчество
<input checked="" type="checkbox"/>	12 Код клиента
<input checked="" type="checkbox"/>	ab buthDay

—>

Фильтрация	
Столбцы присоединяемого набора данных	
<input checked="" type="checkbox"/>	ab Первые 3 буквы
<input checked="" type="checkbox"/>	ab Месяц строковый
<input checked="" type="checkbox"/>	12 Номер месяца

Разбор даты рождения

```

Предпросмотр... | AND | OR | NOT | XOR | = | < | > | <= | >= | 9.0 | " " | 31 | FALSE | TRUE

If(
    RegExMatchedSubExp("^\s*([0-3]?[0-9])[- /.]([0-1]?[0-9])[- /.](19|20)?[0-9]{2}\s*$", Trim(Data_rozhdeniya), 1) <> "",
    RegExMatchedSubExp("^\s*([0-3]?[0-9])[- /.]([0-1]?[0-9])[- /.](19|20)?[0-9]{2}\s*$", Trim(Data_rozhdeniya), 1),
    If(
        RegExMatchedSubExp("^\s*([0-1]?[0-9])[- /.]([0-3]?[0-9])[- /.](19|20)?[0-9]{2}\s*$", Trim(Data_rozhdeniya), 2) <> "",
        RegExMatchedSubExp("^\s*([0-1]?[0-9])[- /.]([0-3]?[0-9])[- /.](19|20)?[0-9]{2}\s*$", Trim(Data_rozhdeniya), 2),
        RegExMatchedSubExp("^\s*((19|20)?[0-9]{2})[- /.]([0-1]?[0-9])[- /.]([0-3]?[0-9])\s*$", Trim(Data_rozhdeniya), 4)
    )
)

```

Объединение

Объединение

№	Главная таблица	<input checked="" type="checkbox"/>	Присоединяемая таблица
1	ab birthDay	<input checked="" type="checkbox"/>	ab buthDay
2	ab birthMonth	<input checked="" type="checkbox"/>	ab burthMonth
3	ab birthYear	<input checked="" type="checkbox"/>	ab birthYear
4	ab Дата рождения	<input checked="" type="checkbox"/>	ab Дата рождения
5	ab Имя	<input checked="" type="checkbox"/>	ab Имя
6	ab Фамилия	<input checked="" type="checkbox"/>	ab Фамилия
7	ab Отчество	<input checked="" type="checkbox"/>	ab Отчество
8	12 Код клиента	<input checked="" type="checkbox"/>	12 Код клиента

Использовать префиксы

Префикс имени	Union
Префикс метки	Объединение

Поделим на разобранные и некорректные

Фильтрация данных

Состояние входа Не активировано

ab birthDay в списке

Формат дня и года

birthYear > Year(Now())

Фильтр строк

Фильтрация данных

Состояние входа Не активировано

ab birthYear2 не в списке

Дата в формате

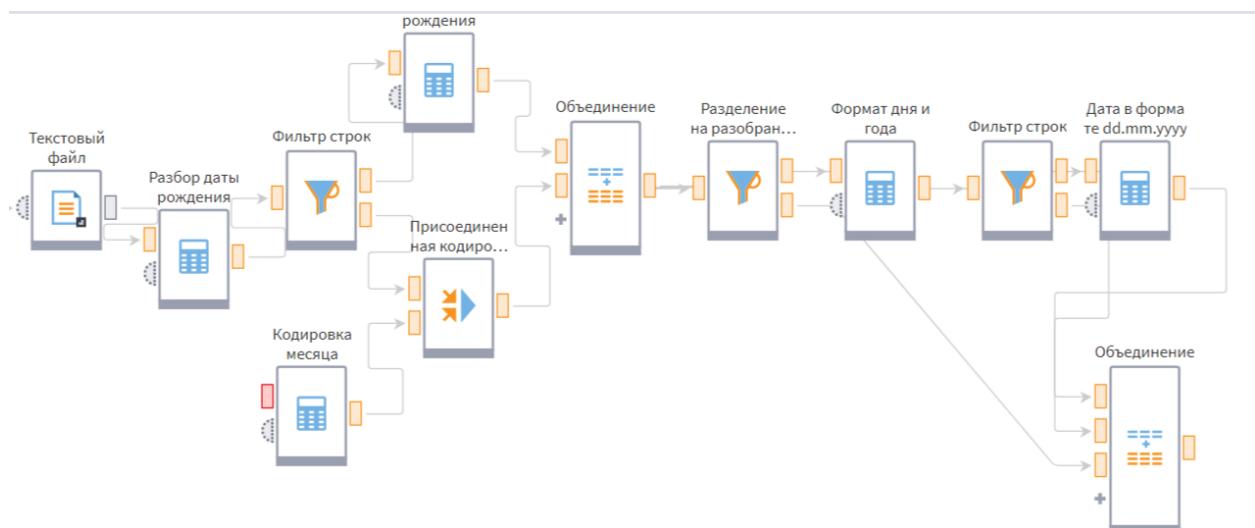
```
Concat(
    birthDay, ".",
    If(
        RegExMatch("^[0-9]{1,2}$", Trim(birthMonth)) = 1,
        Trim(birthMonth),
        IF(Nomer_mesyatsa > 10, Str(Nomer_mesyatsa), CONCAT("0", Str(Nomer_mesyatsa)))
    ),
    ".",
    birthYear
)
```

Объединение

Объединение

№	Главная таблица		Присоединяемая таблица		Присоединяемая таблица
1	ab strResultDate		Не выбрано		Не выбрано
2	31 tesultDate		Не выбрано		Не выбрано
3	ab birthYear		Не выбрано		Не выбрано
4	ab birthMonth		Не выбрано		Не выбрано
5	ab birthDay		Не выбрано		Не выбрано
6	12 Номер месяца		Не выбрано		Не выбрано
7	ab birthYear2		Не выбрано		Не выбрано
8	0/1 invalidYear		Не выбрано		Не выбрано
9	ab birthDay2		Не выбрано		Не выбрано
10	ab Дата рождения		Не выбрано		Не выбрано
11	ab Имя		Не выбрано		Не выбрано
12	ab Фамилия		Не выбрано		Не выбрано
13	ab Отчество		Не выбрано		Не выбрано
14	12 Код клиента		Не выбрано		Не выбрано

Результат



объединение • Быстрый просмотр

дной набор данных

	ab strResultDate	31 resultDate	ab birthY
1	05.10.1986	05.10.1986, 00:00	1986
2	10.05.1986	10.05.1986, 00:00	1986
3	05.05.1986	05.05.1986, 00:00	1986
4	25.05.1986	25.05.1986, 00:00	1986
5	04.05.1986	04.05.1986, 00:00	1986
6	05.05.1952	05.05.1952, 00:00	1952
7	05.01.1986	05.01.1986, 00:00	1986
8	05.10.1906	05.10.1906, 00:00	1906
9	05.10.1986	05.10.1986, 00:00	1986
10	05.10.1986	05.10.1986, 00:00	1986
11	05.10.1987	05.10.1987, 00:00	1987
12	09.05.1986	09.05.1986, 00:00	1986
13	13.05.1986	13.05.1986, 00:00	1986
14	25.05.1987	25.05.1987, 00:00	1987
15	04.05.1987	04.05.1987, 00:00	1987
16	05.05.1953	05.05.1953, 00:00	1953
17	05.01.1987	05.01.1987, 00:00	1987
18	05.10.1966	05.10.1966, 00:00	1966
19	05.10.1946	05.10.1946, 00:00	1946
20	05.10.1987	05.10.1987, 00:00	1987
21	05.10.1988	05.10.1988, 00:00	1988
22	17.05.1986	17.05.1986, 00:00	1986
04	◀		