

Тема 6. Предобработка данных

Синтаксис

Метод `set_axis()` для изменения названий столбцов

```
1 # аргументы - список новых названий столбцов,  
2 # axis со значением columns для изменений в столбцах,  
3 # inplace со значением True для изменения структуры данных  
4 df.set_axis(['a', 'b', 'c'], axis = 'columns', inplace = True)
```

Методы `isnull()` и `isna()` для определения пропущенных значений. В сочетании с методом `sum()` – подсчёт пропущенных значений.

```
1 df.isnull().sum()
```

Метод `fillna()` для заполнения пропущенных значений.

```
1 df = df.fillna(0) # аргумент - значение, на которое будут заменены  
2                 # пропущенные значения
```

Метод `dropna()` для удаления пропущенных значений

```
1 df.dropna() # удаление всех строк, где есть хотя бы одно  
2            # пропущенное значение
```

или

```
1 # аргумент subset - названия столбцов, где нужно искать пропуски  
2 df.dropna(subset = ['a', 'b', 'c'], inplace = True)
```

или

```
1 # аргумент axis со значением 'columns' для удаления столбцов с  
2 # хотя бы одним пропущенным значением  
3 df.dropna(axis = 'columns', inplace = True)
```

Метод `duplicated()` для нахождения дубликатов. В сочетании с методом `sum()` - возвращает количество дубликатов.

```
1 df.duplicated().sum()
```

Метод `drop_duplicates()` для удаления дубликатов. При вызове метода `drop_duplicates()` вместе с повторяющимися строками удаляются их индексы. Поэтому используется с методом `reset_index()`.

```
1 # аргумент drop со значением True,  
2 # чтобы не создавать столбец со старыми значениями индексов  
3 df.drop_duplicates().reset_index(drop = True)
```

Метод `unique()` для просмотра всех уникальных значений в столбце

```
1 df['column'].unique()
```

Метод `replace()` для замены значений в таблице или столбце.

```
1 # первый аргумент - текущее значение
2 # второй аргумент - новое значение
3 df.replace('first_value', 'second_value')
```

Словарь

GIGO (от англ. garbage in — garbage out, буквально «мусор на входе — мусор на выходе») — принцип, который значит, что при неверных входных данных даже правильный алгоритм анализа выдаёт неверные результаты.

Предобработка — процесс подготовки данных для дальнейшего анализа. Заключается она в поиске проблем, которые могут быть в данных, и в устранении этих проблем.

Таблица, удобная для анализа данных:

- в каждом столбце хранятся значения одной переменной;
- каждая строка содержит одно наблюдение, к которому привязаны значения разных переменных.

Названия столбцов:

- Без пробелов в начале, середине и конце.
- Несколько слов разделяются нижним подчеркиваем.
- На одном языке и в одном регистре.
- Отражают в краткой форме, какого рода информация содержится в каждом столбце.

Пропущенные значения бывают разные:

- чаще всего это None или NaN;
- плейсхолдеры (тексты-заполнители) какого-нибудь общепринятого стандарта, иногда неизвестного вам, но которого придерживаются составители. Чаще всего это **n/a**, **na**, **NA**, и **N.N.** либо **NN**.
- произвольное значение, которое по договорённости между собой используют создатели исходной таблицы данных

Пропущенные значения можно как удалять, так и заполнять на основе известных данных.

- Плюс удаления данных в том, что это простой процесс. Также можно быть уверенным, что те данные, которые остались, хорошие и отвечают всем требованиям. Потенциальные минусы: потеря важной информации и снижение точности.
- Заполнение позволяет сохранить наибольшее количество данных. Очевидный минус — могут получиться плохие результаты на основе уже существующих данных.

Дубликаты (дублированные записи) могут быть следующего вида:

- две и более одинаковых строки с идентичной информацией. Большое количество повторов раздувает размер таблицы, а значит, увеличивает время обработки данных;
- одинаковые по смыслу категории с разными названиями, например, «Политика» и «Политическая ситуация». Замаскированные повторы — источник серьёзных и с трудом обнаруживаемых ошибок в анализе.