

Тема 5. Pandas для анализа данных

Синтаксис

Вызов библиотеки pandas

```
1 import pandas
2 import pandas as pd
```

Конструктор DataFrame() для создания таблицы

```
1 pd.DataFrame(data = data, columns = columns)
2 # аргумент data - список с данными,
3 # аргумент columns - список с названиями столбцов
```

Метод read_csv() для чтения файлов формата CSV

```
1 df = pd.read_csv('путь к файлу')
```

Метод head() для вывода первых строк таблицы

```
1 df.head() # первые 5 строк
2 df.head(10) # первые 10 строк
```

Метод tail() для вывода последних строк таблицы

```
1 df.tail() # последние 5 строк
2 df.tail(15) # последние 15 строк
```

Атрибут columns для вывода названий столбцов

```
1 df.columns
```

Атрибут shape для вывода размера таблицы

```
1 df.shape
```

Атрибут dtypes для получения информации о типах данных в таблице

```
1 df.dtypes
```

Метод info() для просмотра сводной информации о таблице

```
1 df.info()
```

Атрибут `loc[строка, столбец]` даёт доступ к элементу в **DataFrame** по строке и столбцу.

```
1 df.loc[:, 'column']
```

Вид	Реализация
Одна ячейка	<code>.loc[7, 'column']</code>
Один столбец	<code>.loc[:, 'column']</code>
Несколько столбцов	<code>.loc[:, ['column_1', 'column_4']]</code>
Несколько столбцов подряд (срез)	<code>.loc[:, 'column_5': 'column_8']</code>
Одна строка	<code>.loc[1]</code>
Все строки, начиная с заданной	<code>.loc[1:]</code>
Все строки до заданной	<code>.loc[:3]</code>
Несколько строк подряд (срез)	<code>.loc[2:5]</code>

Логическая индексация для получения элементов по определенному условию.

Вид	Реализация	Сокращенная запись
Все строки, удовлетворяющие условию	<code>'df.loc[df.loc[:, 'column'] == 'X']'</code>	<code>'df[df['column'] == 'X']'</code>
Столбец, удовлетворяющий условию	<code>'df.loc[df.loc[:, 'column'] == 'X']['column']'</code>	<code>'df[df['column'] == 'X']['column']'</code>
Применение метода	<code>'df.loc[df.loc[:, 'column'] == 'X']['column'].count()'</code>	<code>'df[df['column'] == 'X']['column'].count()'</code>

Индексация в Series

Словарь

Библиотека – это набор готовых методов для решения распространенных задач.

CSV (от англ. Comma-Separated Values, «значения, разделённые запятой») – формат файла. Каждая строка представляет собой одну строку таблицы, где данные разделены запятыми. В первой строке собраны заголовки столбцов (если они есть).

Кортеж – одномерная неизменяемая последовательность данных. Она похожа на список, её тоже можно сохранять в переменной.

DataFrame – это двумерная структура данных **Pandas**, где у каждого элемента есть два индекса: по строке и по столбцу.

- Каждая строка – это одно наблюдение, запись об объекте исследования. А столбцы – признаки этого объекта.
- `DataFrame()` – это конструктор библиотеки **Pandas**, который используется для создания **DataFrame**. Перед именем конструктора стоит обращение к переменной, в которой библиотека хранится: `pd.DataFrame()`.
- У **DataFrame** есть неотъемлемые свойства, значения которых можно запросить. Они называются атрибуты. Например, это размер таблицы `df.shape` или количество столбцов `df.columns`.
- К каждой ячейке с данными в **DataFrame** можно обратиться по её индексу и названию столбца. Этот процесс называется **индексация** и для **DataFrame** его проводят разными способами.

Series — одномерная структура данных Pandas, её элементы можно получить по индексу. Каждый индекс представляет собой номер отдельного наблюдения, и поэтому несколько различных Series вместе составляют DataFrame.

- В Series хранятся данные одного типа.
- У Series есть имя (Name), информация о количестве данных в столбце (Length) и тип данных, которые хранятся в ней (dtype).
- Индексация в Series аналогична индексации элементов столбца в DataFrame.