

# Шпаргалка по Pandas / Pandas Cheatsheet

Краткий справочник по основным операциям

## Содержание

1	Загрузка и Просмотр Данных	1
1.1	А Чтение и Запись	1
1.2	В Первичный Осмотр DataFrame	1
2	Выборка Данных	1
2.1	А Базовый Выбор Столбцов	1
2.2	В Доступ по Меткам и Позициям: .loc vs .iloc	1
3	Фильтрация (Boolean Indexing)	1
3.1	А Фильтрация по Условиям	1
4	Сортировка	2
4.1	А Сортировка DataFrame	2
5	Работа с Дубликатами	2
5.1	А Обнаружение и Удаление Дубликатов	2
6	Применение Функций и Агрегация	2
6.1	А Применение Функций к Элементам, Строкам, Столбцам	2
6.2	В Встроенные Агрегирующие Функции	2
6.3	С Работа со Строками (.str)	2
6.4	D Работа с Датами/Временем (.dt)	2
7	Группировка и Сводные Таблицы	2
7.1	А GroupBy: Разделяй-Применяй-Объединяй	2
7.2	В Сводные Таблицы и Таблицы Сопряженности	3
8	Работа с Пропусками (NaN)	3
8.1	А Обработка Пропущенных Значений (NaN)	3
9	Объединение DataFrames	3
9.1	А Комбинирование DataFrames	3
10	Изменение Формы (Reshaping)	3
10.1	А Изменение Структуры: stack, unstack, melt	3
11	Временные Ряды (Time Series)	3
11.1	А Основные Операции с Временными Рядами	3
12	Преобразование Типов Данных	3
12.1	А Изменение Типов Данных Столбцов	3
13	Базовая Визуализация	4
13.1	А Быстрая Визуализация с .plot()	4

## 1 Загрузка и Просмотр Данных

### 1.1. А Чтение и Запись

#### Основные операции ввода/вывода

Чтение данных из файлов различных форматов и запись DataFrame в файлы.

#### Чтение / Запись

```
1 df_csv = pd.read_csv('file.csv', sep=',')
2 df_excel = pd.read_excel('file.xlsx')
3 df.to_csv('output.csv', index=False)
4 df.to_excel('output.xlsx', index=False)
```

### 1.2. В Первичный Осмотр DataFrame

#### Базовые атрибуты и методы

Команды для получения общего представления о структуре, типах данных и статистиках DataFrame ('df').

#### Осмотр DataFrame ('df')

```
1 df.head()
2 df.tail(3)
3 df.shape
4 df.info()
5 df.describe()
6 df.describe(include='object')
7 df.columns
8 df.dtypes
```

## 2 Выборка Данных

### 2.1. А Базовый Выбор Столбцов

#### Доступ к столбцам по имени

Выбор одного или нескольких столбцов DataFrame по их названию.

#### Выбор столбцов

```
1 s = df['col_name']
2 df_subset = df[['col1', 'col2']]
```

### 2.2. В Доступ по Меткам и Позициям: .loc vs .iloc

#### Ключевое различие

Два основных метода для индексации и выборки данных в Pandas:

- .loc: Выборка по **меткам** индекса и **названиям** столбцов. Правая граница среза **включается**.
- .iloc: Выборка по **целочисленным позициям** (индексам строк и столбцов, начиная с 0). Правая граница среза **не включается** (как в Python).

#### Примеры использования .loc и .iloc

#### Примеры .loc и .iloc

```
1 df.loc['label']
2 df.loc['start':'end']
3 df.loc[:, 'col_name']
4 df.loc['label', 'col_name']
5 df.loc[['l1', 'l3'], ['c1', 'c2']]
6 df.iloc[0]
7 df.iloc[0:5]
8 df.iloc[:, 0]
9 df.iloc[0, 0]
10 df.iloc[[0, 2], [0, 1]]
```

## 3 Фильтрация (Boolean Indexing)

### 3.1. А Фильтрация по Условиям

## Выборка строк на основе логических условий

Используйте логические операторы & (И), | (ИЛИ), ~ (НЕ) для комбинирования условий. Условия необходимо оборачивать в круглые скобки (). Также полезны методы `isin()` и `between()`.

### Примеры Булевой Индексации

```
1 df[df['score'] > 50]
2 df[(df['score'] > 50) & (df['attempts'] < 3)]
3 df[(df['city'] == 'Moscow') | (df['city'] == 'London')]
4 df[df['city'].isin(['London', 'Paris', 'Tokyo'])]
5 df[df['age'].between(18, 30)]
6 df.loc[df['score'] > 90, ['name', 'score']]
7 df[~(df['city'] == 'Moscow')]
```

## 4 Сортировка

### 4.1. А Сортировка DataFrame

#### Упорядочивание данных

Сортировка строк DataFrame по значениям в одном или нескольких столбцах (`sort_values`) или по меткам индекса (`sort_index`).

#### Примеры сортировки

```
1 df.sort_values(by='score')
2 df.sort_values(by='score', ascending=False)
3 df.sort_values(by=['city', 'score'],
4 ascending=[True, False])
5 df.sort_index()
6 df.sort_index(ascending=False)
```

## 5 Работа с Дубликатами

### 5.1. А Обнаружение и Удаление Дубликатов

#### Идентификация и удаление повторяющихся строк

Методы для проверки наличия дубликатов (`uplicated`) и их удаления (`drop_duplicates`).

#### Примеры работы с дубликатами

```
1 df.duplicated()
2 df.duplicated(subset=['user_id', 'timestamp'])
3 df.drop_duplicates()
4 df.drop_duplicates(keep='last')
5 df.drop_duplicates(subset=['user_id'],
6 keep='first')
```

## 6 Применение Функций и Агрегация

### 6.1. А Применение Функций к Элементам, Строкам, Столбцам

#### Поэлементные и построчные/постолбцовые операции

Применение встроенных или пользовательских функций:

- **Поэлементные операции:** Векторизованные операции NumPy или базовые арифметические операции.
- `map()`: Применение функции или словаря к каждому элементу **Series**.
- `apply()`: Применение функции вдоль оси DataFrame (`axis=0` к столбцам, `axis=1` к строкам).
- `applymap()`: Применение функции к каждому элементу **DataFrame**. (Менее рекомендуется, часто есть векторизованные альтернативы).

#### Примеры применения функций

```
1 df['new_col'] = df['col_A'] * 10
2 df['col_B_log'] = np.log(df['col_B'])
3 df['cat_code'] = df['category'].map({'A':1, 'B':2, 'C':3})
4 df[['col_A', 'col_B']].apply(np.sum, axis=0)
5 df['row_sum'] = df[['col_A', 'col_B']].apply(np.sum, axis=1)
6 df['col_C_processed'] = df['col_C'].apply(lambda
7 x: x**2 if x > 0 else 0)
```

### 6.2. В Встроенные Агрегирующие Функции

#### Расчет сводных статистик

Быстрый расчет основных статистик для Series или столбцов DataFrame. Игнорируют NaN по умолчанию.

#### Примеры агрегаций

```
1 df['score'].mean()
2 df['score'].median()
3 df['score'].sum()
4 df['score'].min()
5 df['score'].max()
6 df['score'].std()
7 df['score'].var()
8 df['score'].count()
9 df['score'].nunique()
10 df['category'].unique()
11 df['category'].value_counts()
12 df.mean()
```

### 6.3. С Работа со Строками (.str)

#### Векторизованные строковые операции

Аксессор `.str` предоставляет доступ к множеству методов для работы со строками в Series.

#### Примеры .str методов

```
1 df['name'].str.lower()
2 df['name'].str.upper()
3 df['name'].str.title()
4 df['address'].str.contains('Street')
5 df['email'].str.startswith('info@')
6 df['code'].str.endswith('Z')
7 df['product_id'].str.isdigit()
8 df['city'].str.replace(' ', '-')
9 df['full_name'].str.split(' ')
10 df['comment'].str.len()
```

### 6.4. D Работа с Датами/Временем (.dt)

#### Доступ к компонентам даты/времени

Аксессор `.dt` предоставляет доступ к компонентам даты и времени для Series с типом данных `'datetime64[ns]'`. Предварительно столбец нужно преобразовать.

#### Примеры .dt атрибутов и методов

```
1 df['date_col'].dt.year
2 df['date_col'].dt.month
3 df['date_col'].dt.day
4 df['date_col'].dt.hour
5 df['date_col'].dt.minute
6 df['date_col'].dt.second
7 df['date_col'].dt.dayofweek
8 df['date_col'].dt.dayofyear
9 df['date_col'].dt.quarter
10 df['date_col'].dt.month_name()
11 df['date_col'].dt.day_name()
12 df['date_col'].dt.strftime('%Y-%m-%d')
```

## 7 Группировка и Сводные Таблицы

### 7.1. А GroupBy: Разделяй-Применяй-Объединяй

#### Агрегация данных по группам

Механизм `groupby()` позволяет:

1. **Разделить** данные на группы на основе значений в одном или нескольких столбцах.
2. **Применить** агрегирующую функцию (`sum`, `mean`, `count`, etc.) или трансформацию к каждой группе независимо.

3. **Объединить** результаты в новую структуру данных (Series или DataFrame).

#### Примеры GroupBy

```
1 grouped = df.groupby('category')
2 grouped_multi = df.groupby(['category', 'status'])
3 grouped['value'].sum()
4 grouped['value'].mean()
5 grouped['value'].agg(['mean', 'std', 'count'])
6 grouped_multi.agg(
7     total_value=('value', 'sum'),
8     avg_score=('score', 'mean'),
9     max_score=('score', 'max'),
10    unique_ids=('id', 'nunique')
11 )
12 grouped.size()
```

## 7.2. В Сводные Таблицы и Таблицы Сопряженности

#### Создание агрегированных таблиц

`pivot_table`: Создает сводную таблицу в стиле Excel. Позволяет агрегировать данные по двум или более категориальным переменным. ● `crosstab`: Вычисляет таблицу сопряженности (частот) для двух или более факторов.

#### Pivot Table и CrossTab

```
1 pd.pivot_table(df,
2     values='score',
3     index='category',
4     columns='status',
5     aggfunc=np.mean,
6     fill_value=0)
7 pd.crosstab(df['category'], df['status'])
```

## 8 Работа с Пропусками (NaN)

### 8.1. А Обработка Пропущенных Значений (NaN)

#### Обнаружение, удаление и заполнение NaN

Стандартные подходы к работе с пропущенными данными ('NaN', 'None', 'NaT').

#### Операции с NaN

```
1 df.isnull()
2 df.isna()
3 df.notnull()
4 df.isnull().sum()
5 df.dropna()
6 df.dropna(axis=1)
7 df.dropna(subset=['col_A', 'col_B'])
8 df.dropna(thresh=2)
9 df.fillna(0)
10 df['col_A'].fillna('Unknown')
11 df['col_B'].fillna(df['col_B'].mean())
12 df.fillna(method='ffill')
13 df.fillna(method='bfill')
```

## 9 Объединение DataFrames

### 9.1. А Комбинирование DataFrames

#### Основные методы объединения

- `pd.concat()`: "Склеивание" объектов Pandas вдоль оси (строк `axis=0` или столбцов `axis=1`). Полезно для добавления строк/столбцов.
- `pd.merge()`: Объединение в стиле SQL баз данных. Соединяет строки из двух DataFrame на основе одного или нескольких общих столбцов (**ключей**) или индексов.
- `df.join()`: Удобный метод для объединения по **индексу** (или по ключу в одном DataFrame и индексу в другом). Часто используется как обертка над `merge`.

#### Примеры объединения

Предполагаем наличие 'df1' и 'df2'.

#### Merge, Concat, Join

```
1 combined_rows = pd.concat([df1, df2], axis=0,
2     ignore_index=True)
3 combined_cols = pd.concat([df1, df2], axis=1)
4 merged_inner = pd.merge(df1, df2, on='key',
5     how='inner')
6 merged_left = pd.merge(df1, df2, on='key',
7     how='left')
8 joined_index = df1.join(df2, how='inner',
9     lsuffix='_df1', rsuffix='_df2')
```

## 10 Изменение Формы (Reshaping)

### 10.1. А Изменение Структуры: stack, unstack, melt

#### Преобразование между "широким" и "длинным" форматами

- `stack()`: "Укладывает" столбцы DataFrame в индекс, создавая Series (или DataFrame с MultiIndex). Переход от "широкого" к "длинному" формату.
- `unstack()`: Обратная операция к `stack()`. "Поднимает" уровень индекса в столбцы. Переход от "длинного" к "широкому" формату.
- `melt()`: "Расплавляет" DataFrame, преобразуя столбцы в строки. Полезно для создания "длинного" формата данных из "широкого".

#### Примеры Reshaping

```
1 df_melted = pd.melt(df,
2     id_vars=['id', 'name'],
3     value_vars=['Score_Q1',
4     'Score_Q2'],
5     var_name='Quarter',
6     value_name='Score')
```

## 11 Временные Ряды (Time Series)

### 11.1. А Основные Операции с Временными Рядами

#### Работа с датами и временем как индексом

Pandas предоставляет мощные инструменты для работы с временными рядами, особенно когда индекс DataFrame имеет тип 'DatetimeIndex'.

#### Операции с Time Series

```
1 date_rng = pd.date_range(start='2024-01-01',
2     end='2024-01-10', freq='D')
3 time_rng = pd.date_range('2024-01-01', periods=5,
4     freq='H')
5 df['previous_value'] = df['value'].shift(1)
6 df['change'] = df['value'] - df['value'].shift(1)
```

Доступ к компонентам даты/времени через `.dt` (см. VI.D) работает и для 'DatetimeIndex'.

## 12 Преобразование Типов Данных

### 12.1. А Изменение Типов Данных Столбцов

## Явное приведение типов

Преобразование столбцов к нужным типам данных для корректной обработки или экономии памяти.

### Приведение типов

```
1 print(df.dtypes)
2 df['numeric_col'] =
  pd.to_numeric(df['numeric_col'], errors='coerce')
3 df['string_col'] = df['string_col'].astype(str)
4 df['category_col'] =
  df['category_col'].astype('category')
5 print(df.dtypes)
```

## 13 Базовая Визуализация

### 13.1. А Быстрая Визуализация с .plot()

## Встроенные методы для графиков

Pandas интегрируется с Matplotlib, позволяя быстро строить базовые графики прямо из DataFrame или Series с помощью метода `.plot()`. Необходим импорт `'matplotlib.pyplot'`.

### Примеры `.plot()`

```
1 df['value_ts'].plot(kind='line', title='Trend Over
  Time', figsize=(10, 4))
2 df['numeric_col'].plot(kind='hist', bins=30,
  title='Distribution')
3 df.plot(kind='scatter', x='col_A', y='col_B',
  title='Scatter Plot A vs B')
4 df['numeric_col'].plot(kind='kde', title='Density
  Plot')
```