

Шпаргалка по Pandas / Pandas Cheatsheet (XeLaTeX)

Краткий справочник по основным операциям

April 2, 2025

Contents

- 1 Загрузка и Просмотр
- 2 Выборка: .loc vs .iloc
- 3 Фильтрация (Boolean Indexing)
- 4 Сортировка
- 5 Удаление Дубликатов
- 6 Функции и Агрегация
- 7 Группировка и Сводные Таблицы
- 8 Работа с Пропусками (NaN)
- 9 Объединение DataFrames
- 10 Изменение Формы (Reshaping)
- 11 Временные Ряды (Time Series)
- 12 Преобразование Типов
- 13 Базовая Визуализация

Необходимые импорты

Стандартные импорты для работы с Pandas и NumPy.

Импорты

```
1 import pandas as pd
2 import numpy as np
3 # df = pd.DataFrame(...) # Предполагаемый DataFrame
4 # s = pd.Series(...) # Предполагаемый Series
```

1 Загрузка и Просмотр

Чтение и Запись

Основные операции ввода/вывода данных.

Чтение / Запись

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

Первичный Осмотр

Базовые команды для понимания структуры данных.

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: .loc vs .iloc

Доступ по Меткам и Позициям

- .loc: по меткам индекса/названиям столбцов (включая правую границу среза).
- .iloc: по целочисленным позициям (НЕ включая правую границу среза).

Примеры .loc и .iloc

```
1 # Выбор столбца(ов)
2 s = df['col_name']
3 df_subset = df[['col1', 'col2']]
4
5 # --- .loc (по МЕТКАМ) ---
6 df.loc['label'] # Строка
7 df.loc['start':'end'] # Срез строк (включительно)
8 df.loc[:, 'col_name'] # Столбец
9 df.loc['label', 'col_name'] # Значение
10 df.loc[['l1', 'l3'], ['c1', 'c2']] # По спискам
11
12 # --- .iloc (по ПОЗИЦИЯМ) ---
13 df.iloc[0] # Первая строка
14 df.iloc[0:5] # Первые 5 строк (не вкл. 5)
15 df.iloc[:, 0] # Первый столбец
16 df.iloc[0, 0] # Значение [0, 0]
17 df.iloc[[0, 2], [0, 1]] # По спискам позиций
```

3 Фильтрация (Boolean Indexing)

Выборка по Условиям

Используйте & (И), | (ИЛИ), ~ (НЕ), isin(), between(). Оборачивайте условия в ().

Примеры фильтрации

```
1 df[df['score'] > 50]
2 df[(df['score'] > 50) & (df['attempts'] < 3)]
3 df[df['city'].isin(['London', 'Paris'])]
4 df[df['age'].between(18, 30)]
5 # Явное использование .loc
6 df.loc[df['score'] > 90, ['name', 'score']]
```

4 Сортировка

Упорядочивание Данных

Сортировка по значениям (sort_values) или индексу (sort_index).

Примеры сортировки

```
1 # По значениям
2 df.sort_values(by='score') # Возрастание
3 df.sort_values(by='score', ascending=False) #
  Убывание
4 df.sort_values(by=['city', 'score'],
  ascending=[True, False])
5 # По индексу
6 df.sort_index(ascending=False)
7 # На месте (inplace=True)
8 # df.sort_values(by='score', inplace=True)
```

5 Удаление Дубликатов

Работа с Дубликатами

Обнаружение (duplicated) и удаление (drop_duplicates).

Примеры работы с дубликатами

```
1 df.duplicated() # Проверка (bool Series)
2 df.duplicated(subset=['col1', 'col2'])
3 df.drop_duplicates() # Удалить (оставить первое)
4 df.drop_duplicates(keep='last') # Оставить
  последнее
5 df.drop_duplicates(subset=['user_id'],
  keep='first')
6 # На месте (inplace=True)
7 # df.drop_duplicates(inplace=True)
```

6 Функции и Агрегация

Применение Функций

apply (строки/столбцы), map (Series), applymap (DataFrame). Встроенные агрегации.

Применение функций

```
1 df['new_col'] = df['col_A'] * 10
2 df[['col_A', 'col_B']].apply(np.sum, axis=0)
3 df['cat_code'] = df['cat_name'].map({'A':1, 'B':2})
4 # numeric_df.applymap(lambda x: x + 1)
```

Встроенные Агрегации

Расчет базовых статистик для Series или столбцов DataFrame.

Агрегирующие функции

```
1 df['score'].mean()
2 df['score'].sum()
3 df['score'].min()
4 df['score'].max()
5 df['score'].count() # Не-NA значения
6 df['category'].nunique() # Уникальные значения
```

Работа со Строками и Датами

Использование аксессоров .str и .dt.

Аксессоры .str и .dt

```
1 # Строки (.str)
2 df['name'].str.lower()
3 df['address'].str.contains('Street')
4 df['city'].str.replace(' ', '_')
5 # Даты (.dt, столбец должен быть datetime)
6 # df['date_col'] = pd.to_datetime(df['date_col'])
7 df['date_col'].dt.year
8 df['date_col'].dt.month_name()
9 df['date_col'].dt.dayofweek # Пн=0, Вс=6
```

7 Группировка и Сводные Таблицы

GroupBy (Разделяй-Применяй-Объединяй)

Группировка данных с последующей агрегацией.

Примеры GroupBy

```
1 grouped = df.groupby('category')
2 grouped_multi = df.groupby(['cat', 'status'])
3 grouped['value'].sum()
4 grouped['value'].agg(['mean', 'std', 'count'])
5 grouped_multi.agg({
6     'value': 'sum',
7     'score': ['mean', 'max'],
8     'id': 'nunique'})
9 grouped.size() # Размер групп
10 # df['g_mean'] =
    df.groupby('cat')['val'].transform('mean')
```

Сводные Таблицы

Использование pivot_table и crosstab для агрегации и анализа частот.

Pivot Table и CrossTab

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

8 Работа с Пропусками (NaN)

Обработка NaN

Обнаружение, удаление и заполнение пропущенных данных.

Операции с NaN

```
1 df.isnull().sum() # Количество NaN
2 df.dropna() # Удалить строки с NaN
3 df.dropna(axis=1) # Удалить столбцы с NaN
4 df.dropna(subset=['col_A']) # По столбцу
5 df.fillna(0) # Заполнить нулем
6 df['col_A'].fillna(df['col_A'].mean()) # Средним
7 df.fillna(method='ffill') # Пред. значением
8 df.fillna(method='bfill') # След. значением
```

9 Объединение DataFrames

Merge, Concat, Join

Комбинирование таблиц. merge (SQL-like), concat (склеивание), join (по индексам).

Примеры объединения

```
1 # Merge (по ключу 'key')
2 pd.merge(df1, df2, on='key', how='inner')
3 # Concat (склеить строки)
4 pd.concat([df1, df2], axis=0)
5 # Concat (склеить столбцы)
6 pd.concat([df1, df2], axis=1)
7 # Join (по индексам)
8 df1.join(df2, how='inner')
```

10 Изменение Формы (Reshaping)

Stack, Unstack, Melt

Преобразование структуры DataFrame.

Reshaping

```
1 # Stack (столбцы -> индекс)
2 stacked = df.stack()
3 # Unstack (индекс -> столбцы)
4 unstacked = stacked.unstack()
5 # Melt ("расплавление" столбцов)
6 pd.melt(df, id_vars=['id'], value_vars=['A', 'B'],
7        var_name='variable', value_name='value')
```

11 Временные Ряды (Time Series)

Работа с Датами/Временем

Преобразование, доступ к компонентам, создание диапазонов, смещение, ресемплинг.

Операции Time Series

```
1 # Преобразование и установка индекса
2 df['time_col'] = pd.to_datetime(df['time_col'])
3 df.set_index('time_col', inplace=True)
4 # Доступ к компонентам (если индекс)
5 # year = df.index.year
6 # Создание диапазона
7 rng = pd.date_range('2024-01-01', periods=5,
8                    freq='D')
9 # Смещение
10 df.shift(1)
11 # Ресемплинг (агрегация)
12 df['value'].resample('M').sum() # По месяцам
13 df['value'].resample('W').mean() # По неделям
```

12 Преобразование Типов

Изменение Типов Данных

Явное приведение столбцов к нужным типам (astype, to_numeric).

Приведение типов

```
1 # В числовой (с заменой ошибок на NaN)
2 df['num'] = pd.to_numeric(df['num'],
3                          errors='coerce')
4 # В строку
5 df['col_str'] = df['col_str'].astype(str)
6 # В категорию (экономия памяти)
7 df['cat'] = df['cat'].astype('category')
8 # В datetime (см. раздел Time Series)
9 df['date'] = pd.to_datetime(df['date'])
10 df.dtypes # Проверка типов
```

13 Базовая Визуализация

Быстрые Графики '.plot()'

Создание простых графиков для EDA (требуется 'matplotlib').

Примеры '.plot()'

```
1 # import matplotlib.pyplot as plt
2 # Линейный график
3 df['value'].plot(kind='line', title='Trend')
4 # Гистограмма
5 df['num_col'].plot(kind='hist', bins=20)
6 # Столбчатая диаграмма
7 df.groupby('cat')['val'].mean().plot(kind='bar')
8 # Диаграмма рассеяния
9 df.plot(kind='scatter', x='col_A', y='col_B')
10 # plt.show() # Показать график
```