

pandas cheat sheet

In data science and analysis, pandas is an indispensable library and a great toolkit for data wrangling. Its ability to handle, clean, transform, and analyze data with remarkable efficiency has cemented its position as a cornerstone for both aspiring and seasoned professionals in the field. Consider pandas your trusted companion, because it has a rich array of functions and syntax that unlock hidden narratives within datasets. From completing basic operations to intricate transformations, this cheat sheet guides you through data manipulation with confidence.

Creating and inspecting DataFrames

At the heart of pandas lies the DataFrame, a powerful data structure organizing data into neat rows and columns, mirroring how we conceptualize information. Think of it like a spreadsheet where each row represents a unique entry where these entries could be a person or a transaction, and each column represents a specific characteristic of those entries like age, name, or purchase amount. pandas offers diverse ways to create DataFrames, accommodating various data sources. These sources can come from dictionaries, lists of lists, and CSV files.

Dictionaries are ideal for structured data, where keys represent column names and values correspond to column data, streamlining DataFrame creation from pre-organized data. This is particularly useful when you have data stored in a dictionary format, perhaps after processing JSON data or extracting information from a database.

```
import pandas as pd
```

```
data = {'Name': ['Alice', 'Bob', 'Charlie'], 'Age': [25, 30, 28]}
```

```
df = pd.DataFrame(data)
```

List of lists are convenient for tabular data. Each inner list represents a row, and an optional columns argument specifies column names, ensuring clarity and interpretability. This method is handy when you have data organized in a list-of-lists format, perhaps after reading data from a text file or scraping a website.

```
data = [['Alice', 25], ['Bob', 30], ['Charlie', 28]]
```

```
df = pd.DataFrame(data, columns=['Name', 'Age'])
```

CSV files give pandas the ability to seamlessly load data from them into DataFrames using the read_csv function. This function handles parsing and structuring. CSV files are a common format for storing and exchanging data, and pandas makes it effortless to import them into your Python environment for analysis.

```
df = pd.read_csv('data.csv')
```

pandas also provides similar functions like read_excel, read_json, and read_sql for importing data from other popular formats, making it adaptable to your specific data sources.

Inspecting a DataFrame

Once you have your DataFrame ready, pandas provides a suite of functions to explore and understand its structure and content:

- *df.head()*: Displays the first few rows (default 5) for a quick overview of the data's nature, including column names, data types, and a sample of values. It's like taking a peek at the beginning of a dataset to get a feel for what it contains.
- *df.tail()*: Similarly, displays the last few rows (default 5), aiding in identifying patterns or anomalies towards the end. This can be helpful in spotting any unexpected values or trends that might require further investigation.
- *df.shape*: Returns a tuple representing the dimensions of the DataFrame (number of rows and columns). This is crucial for understanding the scale of your data and planning your analysis strategy accordingly.
- *df.info()*: Generates a concise summary of the DataFrame, encompassing column names, data types (essential for knowing how to work with each column), and the presence of missing values. This summary provides a quick health check of your data, highlighting potential areas for cleaning or preprocessing.
- *df.describe()*: Computes descriptive statistics for numerical columns, offering insights into data distribution and central tendencies. This includes metrics like count, mean, standard deviation, minimum, maximum, and quartiles, which help you understand the characteristics of your numerical data.

Selecting and filtering data

The true power of pandas lies in its ability to precisely select and filter data, allowing you to focus your analysis on specific subsets of interest. This capability is fundamental for extracting meaningful information from large and complex datasets. pandas provides intuitive syntax for column selection.

Single column selection retrieves a single column by its name using square brackets. The result is a pandas Series, a one-dimensional labeled array capable of holding any data type. Series objects inherit many of the DataFrame's functionalities, enabling further manipulation and analysis on a single column.

```
df['Age']
```

Multiple column selection extracts multiple columns by passing a list of column names within double square brackets. The output is a new DataFrame containing only the selected columns, effectively creating a focused subset of the original data. This is useful when you want to work with a specific group of variables or features.

```
df[['Name', 'Age']]
```

pandas offers two primary methods for row selection, each tailored to specific indexing needs.

`df.loc[]` employs label-based indexing, allowing you to select rows based on their index labels, which can be strings, integers, or even dates. This is particularly helpful when your DataFrame has meaningful index labels that you want to use for selection.

`df.loc[0]` # Select the row with index label 0

`df.iloc[]` leverages integer-based indexing, enabling row selection based on their integer positions (starting from 0). This is useful when you want to select rows based on their order in the DataFrame, regardless of their index labels.

`df.iloc[0]` # Select the first row (position 0)

Filtering data is done with Boolean indexing and serves as a powerful tool for filtering rows based on specific conditions. It involves creating a boolean mask (a Series of True/False values) by applying comparison operators or logical conditions to one or more columns. This mask is then used to select only the rows where the condition evaluates to True, effectively filtering the DataFrame based on your criteria. This allows you to extract subsets of data that meet specific requirements, facilitating targeted analysis.

`df[df['Age'] > 25]` # Filter rows where 'Age' is greater than 25

Querying data is for more complex filtering scenarios involving multiple conditions or intricate logic. The `query` method offers a SQL-like syntax. This enhances readability and expressiveness, allowing you to construct queries that resemble natural language expressions, making your code more intuitive and maintainable. It's particularly beneficial when dealing with intricate filtering criteria that would be cumbersome to express using traditional boolean indexing.

`df.query('Age > 25 and Name == "Bob"')` # Filter rows based on multiple conditions

Handling missing data

Real-world datasets are rarely perfect; missing values, denoted as NaN, are common to find. These missing values can come from various sources, such as data entry errors, sensor malfunctions, or incomplete surveys. pandas provides robust mechanisms for identifying and addressing these missing values, ensuring the integrity of your analysis and preventing potential errors or biases that can skew your results.

`df.isnull()` generates a boolean DataFrame mirroring the original, where True indicates a missing value and False represents a non-missing value. This provides a comprehensive overview of the missing data landscape within your DataFrame, allowing you to quickly assess the extent of what is missing and identify columns or rows that require attention.

`df['Age'].isnull()` applies the same logic to a specific column, returning a boolean series highlighting missing values within that column. This allows you to focus your attention on specific variables and assess for their completeness, which is important for deciding on appropriate imputation strategies.

pandas offers a spectrum of strategies for handling missing data, each tailored to specific scenarios and analytical objectives:

`df.dropna()` removes rows containing any missing values, effectively shrinking the DataFrame but ensuring data completeness. This approach is suitable when missing values are relatively sparse and their removal doesn't significantly impact the representativeness of the data. However, it's important

to exercise caution, as dropping rows can lead to loss of valuable information if not done judiciously.

`df.fillna(0)` replaces missing values with a specified value (e.g., 0), a simple imputation technique suitable for certain cases. This can be useful when you have a reasonable assumption about the likely value of missing data or when you want to avoid discarding rows altogether. However, it's important to consider the potential implications of this imputation on your analysis, as it might introduce bias or distort the underlying patterns in the data.

`df['Age'].fillna(df['Age'].mean(), inplace=True)` is a more sophisticated approach, imputing missing values in a column with its mean. The `inplace=True` argument modifies the DataFrame directly, conserving memory. This method is often preferred when missing values are assumed to be randomly distributed and you want to maintain the overall statistical properties of the data.

It's important to be aware that mean imputation can underestimate variance and potentially mask underlying patterns in the data. Other imputation techniques, such as median imputation or more advanced methods like regression imputation or multiple imputation, might be more appropriate depending on the specific characteristics of your dataset and the nature of what is missing.

pandas is an important tool for anyone navigating the world of data. It offers a comprehensive suite of functions and syntax to streamline data manipulation tasks. By learning these functions, you'll be well-equipped to confront diverse data challenges and unveil valuable insights while you work. Remember, practice is key to solidifying your understanding and proficiency with these tools. Explore real-world datasets and spend time experimenting with all the possibilities that pandas offers.