

PANDAS CHEAT-SHEET

1 INSTALLING AND IMPORTING

Installing	<code>pip install pandas</code>
Importing Convention	<code>import pandas as pd</code>

2 READING AND WRITING DATA

Reading data	<code>df = pd.read_csv(path = 'filename.csv')</code> # can extend for json, excel types too using <code>pd.read_json/pd.read_excel</code> , etc.
	<code>df.to_csv('filename.csv')</code> # can extend for json, excel and such too using <code>df.to_json/df.to_excel</code> , etc.
Writing data	<code>df['col'] > value</code> E.g. <code>df['age'] > 30</code>
	<code>df.loc[(df['col1'] == val1) & (df['col2'] == val2)]</code> E.g. <code>df.loc[(df['month'] == 'January') & (df['year'] == '2022')]</code> # filters out data for january 2022

3 SERIES AND DATAFRAMES

• CREATING A SERIES	
<code>pd.Series(['a', 'b', 'c'])</code>	
• CREATING A DATAFRAME	
Row Oriented	<code>pd.DataFrame([['a', 1], ['b', 2], columns=['name', 'id'])</code>
Column Oriented	<code>pd.DataFrame({'name': ['a', 'b'], 'id':[1, 2]})</code>

4 INFO EXTRACTION

Shape (Return a tuple representing the dimensionality of the DataFrame.)	<code>df.shape</code> #e.g. -(2,3) for 2 rows and 3 columns
Head (first n rows, default 5)	<code>df.head(n)</code>
Tail (last n rows, default 5)	<code>df.tail(n)</code>
Info (return info of all columns)	<code>df.info()</code>
Describe (gives statistical information of data)	<code>df.describe()</code>

5 ACCESSING

• DIRECT ACCESSING COLUMNS AND ROWS, AS WELL AS BOTH TOGETHER	
Accessing a row	<code>df.loc[ei]</code> #(ei here is explicit index) <code>df.iloc[ii]</code> #(ii here is implicit index)
Accessing a column	<code>df['column_name']</code> #for single column <code>df[['col1', 'col2']]</code> #for multiple columns
• SLICING	
Row	<code>df.loc[1:3]</code> (1 and 3 are the explicit indices here) Or <code>df.iloc[2:4]</code> (2 and 4 are the implicit indices here)

Column	<code>df.loc[:, 'a':'b']</code>
Both row and columns	<code>df.loc[1:3, 'a':'b']</code> (1 and 3 are explicit indices here)

• FEATURE EXPLORATION (MASKING, FILTERING)

Masking	<code>df['col'] > value</code> E.g. <code>df['age'] > 30</code>
	<code>df.loc[(df['col1'] == val1) & (df['col2'] == val2)]</code>

Filtering	Filters data based on conditions
	<code>df.loc[(df['month'] == 'January') & (df['year'] == '2022')]</code> # filters out data for january 2022

6 DATFRAME MANIPULATION

• ADDING A NEW ROW/COLUMN

Row	<code>df.loc[explicit_row_num] = ['a', 1]</code> E.g. <code>df.loc[len(df.index)] = ['a', 1]</code> # this will add a row at the end of the dataframe
Column	<code>df['new_col'] = data</code>

• DELETING A NEW ROW/COLUMN

Row	<code>df.drop(labels=None, axis=0)</code> E.g. <code>df.drop(3, axis=0)</code> # Here 3 is the explicit index, axis=0 is for row
Column	<code>df.drop('col_name', axis=1)</code>

• RENAMING A COLUMN

Row	<code>df.index = new_indices</code>
Column	<code>df.rename({'old_name': 'new_name', axis=1})</code>

• DUPLICATES AND DROPPING DUPLICATES

• find duplicate rows

df.duplicated(subset=None, keep='first')	#subset can be used to specify certain column(s) for identifying the duplicates
	# keep determines which duplicates to mark
1. first : marks all duplicates as true except for the first occurrence.	
2. last : marks all duplicates as true except for the last occurrence.	
3. False : marks all duplicates as true	returns a boolean series for each duplicate row marked as true
• drop duplicate values	
<code>df.drop_duplicates(subset=None, keep='first')</code>	# parameters have the same meaning as in df.duplicated, except here it will drop the rows marked duplicate

7 OPERATIONS

• SORTING

<code>df.sort_values(['col1'], ascending=[True])</code>

• BUILT IN OPS

built in ops such as mean, min, max, etc.
e.g. <code>df['col1'].min(), df['col1'].count(), etc.</code>

• APPLY

applies a function along one of the axis of the dataframe
<code>df['col'].apply(function)</code>
e.g. <code>data[['revenue', 'budget']].apply(np.sum, axis=1)</code> #sums values of revenue and budget across each row

Group based filtering	<code>df.groupby('group_col_name').filter(boolean array based on condition)</code>
	E.g. <code>data.groupby('director_name').filter(lambda x: x["budget"].max() >= 100)</code>
	# This filters all rows of those directors whose maximum budget is greater than 100 million

Group based apply	<code>df.groupby('group_col_name').apply(function)</code>
	E.g. <code>def func(x): x["risky"] = x["budget"] - x["revenue"].mean() >= 0 return x</code> <code>data_risky = data.groupby("director_name").apply(func)</code>
	# Finds movies whose budget is higher than its director's average revenue

Pivot	<code>df.pivot(index=['list of columns'], columns='col_name', values='col_name')</code>
	E.g. <code>data_melt.pivot(index=['Date', 'Drug_Name'], columns = 'time', values='reading')</code>



Cut	<code>df['new_cat_column']=pd.cut(df['continuous_col'], bins=bin_values, labels=label_values)</code>
	E.g. <code>data_tidy['temp_cat'] = pd.cut(data_tidy['Temperature'], bins=temp_points, labels=temp_labels)</code>

Shift	<code>df['col'].shift(periods=n, axis=0)</code>
	E.g. <code>df['Marks'].shift(periods = 1, axis = 0)</code>

10 CLEANING OUR DATA	
• NONE AND NAN	
	• "NaN" is for columns with numbers as their values
	• "None" is for columns with non-number entries(e.g. String, object type, etc.)
	• Can check for null values using "isna()"
	• E.g.
	• <code>df.isna()</code> # returns the dataframe with True/False for null values in the respective element's position
	• <code>df.isna().sum()</code> # returns number of null values per column. Can modify with <code>df.isna().sum(axis=1)</code> for each row's null count
	• <code>df.isna().sum().sum()</code> # returns total number of null values
• FILLING NULL VALUES	
	<code>df.fillna(n) # fills null values with value 'n'</code>
• DROPPING NULL VALUES	
	<code>df.dropna(axis=0)</code>
	# Default axis=0, use 1 for columns
	# Drops rows/columns with even a single missing value

||
||
||