

pandas

The everyday pandas moves, from building and cleaning frames to grouping, reshaping, and reading or writing them.

Build a Frame

From dicts, lists, or files into a DataFrame.

PURPOSE	COMMAND	RESULT
From a dict of columns.	<code>pd.DataFrame({"city": cities, "pop": pops})</code>	
One labeled column (Series).	<code>pd.Series([8.4, 3.9, 0.9], name="pop")</code>	
From row records.	<code>pd.DataFrame.from_records(rows)</code>	
Read a CSV file.	<code>pd.read_csv("cities.csv")</code>	
Read Parquet (columnar).	<code>pd.read_parquet("cities.parquet")</code>	
Empty frame with set columns.	<code>pd.DataFrame(columns=["city", "pop"])</code>	

a DataFrame is a dict of aligned, named columns (Series) sharing one index

Look at It



Shape, types, summary, a peek at the data.

PURPOSE	COMMAND	RESULT
First / last n rows.	<code>df.head(5) · df.tail(5)</code>	
Rows and columns count.	<code>df.shape</code>	
Schema, dtypes, non-null.	<code>df.info()</code>	
Per-column dtypes.	<code>df.dtypes</code>	
Numeric summary stats.	<code>df.describe()</code>	
Count distinct values.	<code>df["region"].value_counts()</code>	

look before you transform: shape, dtypes, and where the nulls hide

Pick Rows and Columns

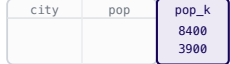





Labels with `.loc`, positions with `.iloc`, masks for filtering.

PURPOSE	COMMAND	RESULT
Select columns by name.	<code>df[["city", "pop"]]</code>	
Select by label (.loc).	<code>df.loc[df["pop"] > 5, "city"]</code>	
Select by integer position.	<code>df.iloc[0:5, 0:2]</code>	
Filter with a boolean mask.	<code>df[df["pop"] > 5]</code>	
Filter with a query string.	<code>df.query("pop > 5 and region == 'E'")</code>	
Keep rows where value in set.	<code>df[df["region"].isin(["E", "W"])]</code>	

`.loc` takes labels and masks; `.iloc` takes integer positions

Make New Columns

Add, rename, retype, and transform values.

PURPOSE	COMMAND	RESULT
Add or replace columns.	<code>df.assign(pop_k=df["pop"] * 1000)</code>	
Rename columns.	<code>df.rename(columns={"pop": "population"})</code>	
Change a column's dtype.	<code>df.astype({"pop": "int64"})</code>	
Elementwise string transform.	<code>df["city"].str.upper()</code>	
Map / apply a function.	<code>df["pop"].apply(round)</code>	
Drop columns.	<code>df.drop(columns=["notes"])</code>	

most calls return a new frame; chain them rather than mutating in place

Fix Missing & Messy Data

Nulls, duplicates, and ordering.

PURPOSE	COMMAND	RESULT						
Count nulls per column.	<code>df.isna().sum()</code>	city 0 pop 2						
Drop rows with any null.	<code>df.dropna()</code>	<table border="1"> <tr><td>NY</td><td>8.4</td></tr> <tr><td>LA</td><td>NaN</td></tr> </table>	NY	8.4	LA	NaN		
NY	8.4							
LA	NaN							
Fill nulls with a value.	<code>df.fillna(0)</code>							
Forward-fill nulls.	<code>df.ffill()</code>							
Remove duplicate rows.	<code>df.drop_duplicates()</code>	<table border="1"> <tr><td>NY E 8.4</td></tr> <tr><td>NY E 8.4</td></tr> </table>	NY E 8.4	NY E 8.4				
NY E 8.4								
NY E 8.4								
Sort by column values.	<code>df.sort_values("pop", ascending=False)</code>	<table border="1"> <thead> <tr><th>city</th><th>pop</th></tr> </thead> <tbody> <tr><td>NY</td><td>8.4</td></tr> <tr><td>LA</td><td>3.9</td></tr> </tbody> </table>	city	pop	NY	8.4	LA	3.9
city	pop							
NY	8.4							
LA	3.9							

`isna().sum()` finds nulls; `dropna` or `fillna/ffill` resolve them

Group and Summarize

Split into groups, aggregate, summarize into tables.

PURPOSE	COMMAND	RESULT									
Group then aggregate.	<code>df.groupby("region")["pop"].sum()</code>	<table border="1"> <tr><td>E</td><td>9.1</td></tr> <tr><td>W</td><td>4.8</td></tr> </table>	E	9.1	W	4.8					
E	9.1										
W	4.8										
Named aggregations.	<code>df.groupby("region").agg(total=("pop", "sum"), n=("pop", "size"))</code>	<table border="1"> <thead> <tr><th>region</th><th>total</th><th>n</th></tr> </thead> <tbody> <tr><td>E</td><td>9.1</td><td>2</td></tr> <tr><td>W</td><td>4.8</td><td>2</td></tr> </tbody> </table>	region	total	n	E	9.1	2	W	4.8	2
region	total	n									
E	9.1	2									
W	4.8	2									
Count rows per group.	<code>df.groupby("region").size()</code>	E : 2 W : 2									
Broadcast a group stat back.	<code>df.groupby("region")["pop"].transform("mean")</code>	<table border="1"> <thead> <tr><th>mean</th></tr> </thead> <tbody> <tr><td>4.55</td></tr> <tr><td>2.40</td></tr> <tr><td>..</td></tr> </tbody> </table> <p>same length as df</p>	mean	4.55	2.40	..					
mean											
4.55											
2.40											
..											
Pivot to a summary table.	<code>df.pivot_table(index="region", values="pop", aggfunc="mean")</code>	<table border="1"> <thead> <tr><th>region</th><th>pop</th></tr> </thead> <tbody> <tr><td>E</td><td>4.55</td></tr> <tr><td>W</td><td>2.40</td></tr> </tbody> </table>	region	pop	E	4.55	W	2.40			
region	pop										
E	4.55										
W	2.40										
Frequency cross-tab.	<code>pd.crosstab(df["region"], df["size"])</code>	<table border="1"> <thead> <tr><th></th><th>S</th><th>L</th></tr> </thead> <tbody> <tr><td>E</td><td>1</td><td>1</td></tr> <tr><td>W</td><td>1</td><td>1</td></tr> </tbody> </table>		S	L	E	1	1	W	1	1
	S	L									
E	1	1									
W	1	1									

split, apply, combine: the group key lands in the result's index

Reshape and Combine

Join frames, stack rows, wide-to-long and back.

PURPOSE	COMMAND	RESULT
Join on a key column.	<code>pd.merge(left, right, on="id", how="left")</code>	
Stack rows (same columns).	<code>pd.concat([df_jan, df_feb])</code>	
Glue columns side by side.	<code>pd.concat([df_a, df_b], axis=1)</code>	
Wide to long (unpivot).	<code>df.melt(id_vars="city", var_name="year", value_name="pop")</code>	
Long to wide (pivot).	<code>df.pivot(index="city", columns="year", values="pop")</code>	
Expand list-valued cells.	<code>df.explode("tags")</code>	

merge joins on a key; melt and pivot are wide/long inverses

Read, Write & Dates

File I/O round-trips and datetime handling.

PURPOSE	COMMAND	RESULT
Write a CSV (no index col).	<code>df.to_csv("out.csv", index=False)</code>	
Write Parquet (typed, fast).	<code>df.to_parquet("out.parquet")</code>	
Parse strings to datetimes.	<code>pd.to_datetime(df["date"])</code>	
Pull out date parts.	<code>df["date"].dt.year .dt.month_name()</code>	
Build a date range index.	<code>pd.date_range("2024-01-01", periods=12, freq="ME")</code>	
Downsample a time series.	<code>df.resample("ME").sum()</code>	

once a column is real datetimes, .dt and resample unlock calendar work