

polars

A quick reference for Polars, the fast DataFrame library with a lazy engine, built around expressions for selecting, filtering, grouping, and joining.

Read & Write Data

Eagerly with `read_*`, or lazily and out-of-core with `scan_*`.

PURPOSE	COMMAND	RESULT
Read a CSV into memory.	<code>pl.read_csv("data.csv")</code>	
Read a Parquet file.	<code>pl.read_parquet("data.parquet")</code>	
Scan lazily (no data read).	<code>pl.scan_parquet("data.parquet")</code>	
Build a frame from a dict.	<code>pl.from_dict({"a": [1, 2], "b": [3, 4]})</code>	
Write to Parquet.	<code>df.write_parquet("out.parquet")</code>	
Stream a query to disk.	<code>lf.sink_parquet("out.parquet")</code>	

`scan_*` defers the read so polars can push filters into the file scan

Inspect a Frame

Look before you leap: shape, schema, a peek, and summary stats.

PURPOSE	COMMAND	RESULT								
Rows and columns.	<code>df.shape</code>	<code>(1000, 5)</code>								
First few rows.	<code>df.head(5)</code>	<table border="1"><thead><tr><th>a</th><th>v</th></tr></thead><tbody><tr><td>1</td><td>10</td></tr><tr><td>2</td><td>20</td></tr><tr><td>..</td><td>..</td></tr></tbody></table>	a	v	1	10	2	20
a	v									
1	10									
2	20									
..	..									
Column names and types.	<code>df.schema</code>	<code>a</code> → Int64 <code>v</code> → Float64 <code>b</code> → String								
Wide, transposed preview.	<code>df.glimpse()</code>	<pre>a<i64> 1, 2, 3 b<str> "x", "y" v<f64> 1.5, 2.0</pre>								
Summary statistics.	<code>df.describe()</code>	<table border="1"><thead><tr><th></th><th>v</th></tr></thead><tbody><tr><td>count</td><td>1000</td></tr><tr><td>mean</td><td>25.0</td></tr><tr><td>max</td><td>99.0</td></tr></tbody></table>		v	count	1000	mean	25.0	max	99.0
	v									
count	1000									
mean	25.0									
max	99.0									
Nulls per column.	<code>df.null_count()</code>	<table border="1"><thead><tr><th>a</th><th>b</th><th>v</th></tr></thead><tbody><tr><td>0</td><td>0</td><td>3</td></tr></tbody></table>	a	b	v	0	0	3		
a	b	v								
0	0	3								

`glimpse` is the wide transposed peek; `describe` and `null_count` flag cleanup work

Select & Add Columns

Expressions are the heart of polars: select to keep, with_columns to add.

PURPOSE	COMMAND	RESULT									
Keep specific columns.	<code>df.select("a", "b")</code>	<table border="1"><thead><tr><th>a</th><th>b</th><th>c</th><th>d</th></tr></thead><tbody><tr><td></td><td></td><td></td><td></td></tr></tbody></table>	a	b	c	d					
a	b	c	d								
Build an expression.	<code>pl.col("v") * 2</code>										
Add / overwrite a column.	<code>df.with_columns((pl.col("v") * 2).alias("v2"))</code>	<table border="1"><thead><tr><th>a</th><th>v</th><th>v2</th></tr></thead><tbody><tr><td></td><td></td><td>20</td></tr><tr><td></td><td></td><td>40</td></tr></tbody></table>	a	v	v2			20			40
a	v	v2									
		20									
		40									
Conditional column.	<code>pl.when(pl.col("v") > 0).then(pl.lit("pos")).otherwise(pl.lit("neg"))</code>										
Change a column's type.	<code>pl.col("a").cast(pl.Float64)</code>										
Select by selector.	<code>df.select(cs.numeric())</code>	<table border="1"><thead><tr><th>i64</th><th>str</th><th>f64</th><th>i64</th></tr></thead><tbody><tr><td></td><td></td><td></td><td></td></tr></tbody></table>	i64	str	f64	i64					
i64	str	f64	i64								

`pl.col("v") * 2` is a recipe, not a value; verbs run it against the frame

Filter & Sort Rows

Pick the rows you want with boolean expressions, then order them.

PURPOSE	COMMAND	RESULT				
Keep rows matching a test.	<code>df.filter(pl.col("v") > 100)</code>	<table border="1"><thead><tr><th>v</th></tr></thead><tbody><tr><td>120</td></tr><tr><td>40</td></tr><tr><td>150</td></tr></tbody></table> 	v	120	40	150
v						
120						
40						
150						
Combine conditions.	<code>df.filter((pl.col("v") > 100) & (pl.col("g") == "a"))</code>					
Keep rows in a set.	<code>df.filter(pl.col("g").is_in(["a", "b"]))</code>					
Sort by a column.	<code>df.sort("v", descending=True)</code>	<table border="1"><thead><tr><th>v</th></tr></thead><tbody><tr><td>40</td></tr><tr><td>30</td></tr><tr><td>20</td></tr></tbody></table>	v	40	30	20
v						
40						
30						
20						
Sort by several keys.	<code>df.sort(["g", "v"], descending=[False, True])</code>					
Drop duplicate rows.	<code>df.unique(subset=["g"])</code>	<table border="1"><tbody><tr><td>a 10</td></tr><tr><td>a 10</td></tr></tbody></table> 	a 10	a 10		
a 10						
a 10						

combine conditions with `&` | `~`, each side wrapped in parentheses

Group & Aggregate

Split by key, aggregate each group, or compute window stats with over.

PURPOSE	COMMAND	RESULT
Sum within each group.	<code>df.groupby("g").agg(pl.col("v").sum())</code>	
Count rows per group.	<code>df.groupby("g").len()</code>	
Several aggregates at once.	<code>df.groupby("g").agg(pl.col("v").mean().alias("avg"), pl.len().alias("n"))</code>	
Group by multiple keys.	<code>df.groupby(["g", "h"]).agg(pl.col("v").sum())</code>	
Stat without collapsing rows.	<code>df.with_columns(pl.col("v").sum().over("g").alias("grp_sum"))</code>	
Time-windowed groups.	<code>df.group_by_dynamic("t", every="1h").agg(pl.col("v").mean())</code>	

`over("g")` computes a group stat but keeps every original row

Join & Concatenate







Combine frames side-by-side on keys, or stack them.

PURPOSE	COMMAND	RESULT
Inner join on a key.	<code>left.join(right, on="k", how="inner")</code>	
Keep all left rows.	<code>left.join(right, on="k", how="left")</code>	
Outer / full join.	<code>left.join(right, on="k", how="full")</code>	
Filtering joins (semi / anti).	<code>left.join(right, on="k", how="semi")</code>	
Stack rows (vertical).	<code>pl.concat([df1, df2])</code>	
Stack columns (horizontal).	<code>pl.concat([df1, df2], how="horizontal")</code>	

joins combine by key; concat just stacks frames by rows or columns

Reshape & Restructure







Pivot wide, unpivot long, explode lists, and rename.

PURPOSE	COMMAND	RESULT
Wide to long (unpivot).	<pre>df.unpivot(index="id", on=["x", "y"])</pre>	
Long to wide (pivot).	<pre>df.pivot("var", index="id", values="val")</pre>	
One list to many rows.	<pre>df.explode("items")</pre>	
Rename columns.	<pre>df.rename({"x": "X"})</pre>	
Drop columns.	<pre>df.drop("tmp")</pre>	
Flip rows and columns.	<pre>df.transpose()</pre>	

unpivot gathers value columns into rows; pivot spreads them back out

Clean: Strings, Dates & Nulls

The everyday tidy-up: namespaced string and date ops, plus null handling.

PURPOSE	COMMAND	RESULT
Transform text.	<pre>pl.col("s").str.to_uppercase()</pre>	
Match a substring.	<pre>pl.col("s").str.contains("ell")</pre>	
Parse text to a date.	<pre>pl.col("d").str.to_date()</pre>	
Extract a date part.	<pre>pl.col("d").dt.year()</pre>	
Fill missing values.	<pre>pl.col("v").fill_null(0)</pre>	
Drop rows with nulls.	<pre>df.drop_nulls()</pre>	

.str and .dt namespaces hold type-specific ops; null is explicit, not NaN