



# Polars Cheat Sheet

[Open in Colab](#)

## General

### Install

pip install polars

### Import

import polars as pl

Polars DataFrame / Series

Polars DataFrame

nrs	names	random	groups
1	"foo"	0.3	"A"
2	"ham"	0.7	"A"
3	"spam"	0.1	"B"
null	"egg"	0.9	"C"
5	null	0.6	"B"

```
df = pl.DataFrame({
    "nrs": [1, 2, 3, None, 5],
    "names": ["foo", "ham", "spam", "egg", None],
    "random": [0.3, 0.7, 0.1, 0.9, 0.6],
    "groups": ["A", "A", "B", "C", "B"],
})
```

### Read CSV

```
df = pl.read_csv("https://j.mp/iris.csv",
    has_header=True)
```

### Read parquet

```
df = pl.read_parquet("path.parquet",
    columns=["select",
    "columns"])
```

## Expressions

Polars Expressions

Polars

```
df \
    .filter(pl.col("nrs") < 4) \
    .group_by("groups") \
    .agg(pl \
    .all() \
    .sum()
)
```

Polars DataFrame



Polars DataFrame

```
df.filter(pl.col("random") > 0.5)
df.filter(
    (pl.col("groups") == "B")
    & (pl.col("random") > 0.5)
)
```

Polars DataFrame

```
# Randomly select fraction of rows.
df.sample(frac=0.5)

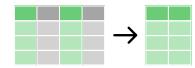
# Randomly select n rows.
df.sample(n=2)
```

Polars DataFrame

```
# Select first n rows
df.head(n=2)

# Select last n rows.
df.tail(n=2)
```

Polars DataFrame



Polars DataFrame

```
df.select(["nrs", "names"])
```

Polars DataFrame

```
df.select(pl.col("^n.*$"))
```

Polars DataFrame



Polars DataFrame

```
df[2:4, :]
```

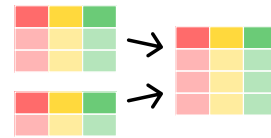
Polars DataFrame

```
df[:, [1, 3]]
```

Polars DataFrame

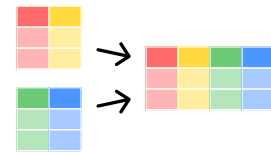
```
df[df["random"] > 0.5, ["names", "groups"]]
```

Polars DataFrame



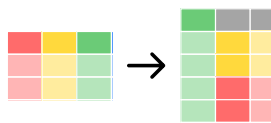
Polars DataFrame

```
pl.concat([df, df2])
```



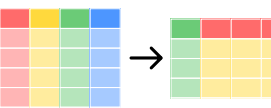
Polars DataFrame

```
pl.concat([df, df3], how="horizontal")
```



Polars DataFrame

```
df.melt(
    id_vars="nrs",
    value_vars=["names", "groups"]
)
```



Polars DataFrame

```
df.pivot(values="nrs", index="groups",
    columns="names")
```

Polars DataFrame

```
# low to high
df.sort("random")

# high to low
df.sort("random", reverse=True)
```

Polars DataFrame

```
df.rename({"nrs": "idx"})
```

Polars DataFrame

```
df.drop(["names", "random"])
```

Polars DataFrame

Polars DataFrame

```
df["groups"].value_counts()
```

Polars DataFrame

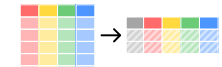
```
len(df)
# or
df.height
```

Polars DataFrame

```
df.shape
```

Polars DataFrame

```
df["groups"].n_unique()
```



Polars DataFrame

```
df.describe()
```

Polars DataFrame

```
df.select(
    [
        # Sum values
        pl.sum("random").alias("sum"),

        # Minimum value
        pl.min("random").alias("min"),

        # Maximum value
        pl.max("random").alias("max"),
        # or
        pl.col("random").max().alias("other_max"),

        # Standard deviation
        pl.std("random").alias("std dev"),

        # Variance
        pl.var("random").alias("variance"),

        # Median
        pl.median("random").alias("median"),

        # Mean
        pl.mean("random").alias("mean"),

        # Quantile
        pl.quantile("random", 0.75) \
        .alias("quantile_0.75"),
        # or
        pl.col("random").quantile(0.75) \
        .alias("other_quantile_0.75"),

        # First value
        pl.first("random").alias("first"),
    ]
)
```



```
df.group_by(by="groups").agg(
[
    # Sum values
    pl.sum("random").alias("sum"),

    # Minimum value
    pl.min("random").alias("min"),

    # Maximum value
    pl.max("random").alias("max"),
    # or
    pl.col("random").max().alias("other_max"),

    # Standard deviation
    pl.std("random").alias("std_dev"),

    # Variance
    pl.var("random").alias("variance"),

    # Median
    pl.median("random").alias("median"),

    # Mean
    pl.mean("random").alias("mean"),
    # Quantile
    pl.quantile("random", 0.75) \
        .alias("quantile_0.75"),
    # or
    pl.col("random").quantile(0.75) \
        .alias("other_quantile_0.75"),
    # First value
    pl.first("random").alias("first"),
]
```



The diagram shows a 4x4 matrix on the left with missing values indicated by diagonal lines. An arrow points to the right, where the same matrix is shown with the missing values filled in, labeled as 'value'.



```
df.select(
    [
        "names",
        "groups",
        pl.col("random").sum().over("names") \
        .alias("sum_by_names"),
        pl.col("random").sum().over("groups") \
        .alias("sum_by_groups"),
    ]
)
```

nrs	names
1	"foo"
2	"ham"
3	"spam"

 $\bowtie$ 

nrs	animals
1	"cheetah"
2	"lion"
6	"tiger"

 $=$ 

nrs	names	animals
1	"foo"	"cheetah"
2	"ham"	"lion"



nrs	names
1	"foo"
2	"ham"
3	"spam"

 $\bowtie$ 

nrs	animals
1	"cheetah"
2	"lion"
6	"tiger"

 $=$ 

nrs	names	animals
1	"foo"	"cheetah"
2	"ham"	"lion"
3	"spam"	<del>null</del>

The diagram shows two input tables being joined. The first table has columns 'nrs' and 'names' with rows (1, 'foo'), (2, 'ham'), and (3, 'spam'). The second table has columns 'nrs' and 'animals' with rows (1, 'cheetah'), (2, 'lion'), and (6, 'tiger'). The result of the join is a table with columns 'nrs' and 'names' containing the row (3, 'spam').

