# Polars Cheat Sheet

Open in Colab

## General

### Install

```
pip install polars
```

### Import

```
import polars as pl
```

## データフレームの作成/読み込み

### データフレームの作成

```
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"],
    }
)
```

| 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" |

### Read CSV

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

### Read parquet

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

## Expressions

Polarsの式は柔軟で構成可能です。 これにより、効率的な処理を実現します。
以下は例:

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

## サブセットの取得 - 行

### 条件に一致する: 条件に一致する行を抽出する

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

### ランダム

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

# Randomly select n rows.
df.sample(n=2)
```
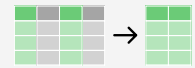
### 先頭・末尾の行を取得する

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

# Select last n rows
df.tail(n=2)
```
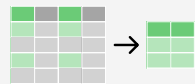
## サブセットの取得 - 列

### 列名で複数の列を選択する

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

### 正規表現で複数の列を選択する

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

## サブセット - 行と列

### 2~4行目を取得

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

### 第1と3列を取得 (先頭は0)

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

### 条件に一致する行と特定の列を取得する

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

## データフレームの結合・整形

### データフレームを縦に結合する

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

### データフレームを横に結合する

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

### 縦持ちに変換する

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

### 横持ちに変換する

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

### 値でソートする

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

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

### 列名を変更する

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

### 列を削除する

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

## データの要約

### 出現回数をカウントする

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

### 行数をカウントする

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

### 行数・列数を取得する

```
df.shape
```

### ユニークな値の数

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

### 基本統計量を取得する

```
df.describe()
```

### 集計

```
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"),
    ]
)
```

□□□ col □□□□□□□□□□□□ GroupBy □□□□□□□□□

```python
df.group_by("groups")
```

□□□□□□□□□□□□□□□□□□□□□□□

```python
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"),
    ]
)
```

□□□ GroupBy □□□

```python
df.group_by(by="groups").agg(
    [
        # Count the number of values in each group
        pl.count("random").alias("size"),

        # Sample one element in each group
        pl.col("names").apply(
            lambda group_df: group_df.sample(1)
        ),
    ]
)
```

---

□□□□□□□□□



□□□□□ null □□□□□□□□□□□□□□

```python
df.drop_nulls()
```



□□□□□□□ null □□□□□□

```python
df.fill_null(42)
```



□□□□□□□□□□ null □□□□□□

```python
df.fill_null(strategy="forward")
```

□□□□□□□□□ "backward", "min", "max", "mean",
"zero" "one" □□□□□□□

□□□□□□ NaN □□□□□□□□□□□□□□

```python
df.fill_nan(42)
```

□□□□□□□□□□□



□□□□□□□□□□□□□□□□□□

```python
df.with_column(
    (pl.col("random") * pl.col("nrs")) \
        .alias("product")
)
```

□□□□□□□□□□□□□□□□□□□□□□

```python
df.with_columns(
    [
        (pl.col("random") * pl.col("nrs")) \
            .alias("product"),
        pl.col("names").str.lengths() \
            .alias("names_lengths"),
    ]
)
```

□□□□□□ 0 □□□□□□□□□□□□□□□□

```python
df.with_row_count()
```

---

□□□□□□□□



□□□□□□□□□□□□□□□

```python
df.select(
    [
        # Rolling maximum value
        pl.col("random") \
            .rolling_max(window_size=2) \
            .alias("rolling_max"),

        # Rolling mean value
        pl.col("random") \
            .rolling_mean(window_size=2) \
            .alias("rolling_mean"),

        # Rolling median value
        pl.col("random") \
            .rolling_median(
                window_size=2, min_periods=2) \
            .alias("rolling_median"),

        # Rolling minimum value
        pl.col("random") \
            .rolling_min(window_size=2) \
            .alias("rolling_min"),

        # Rolling standard deviation
        pl.col("random") \
            .rolling_std(window_size=2) \
            .alias("rolling_std"),

        # Rolling sum values
        pl.col("random") \
            .rolling_sum(window_size=2) \
            .alias("rolling_sum"),

        # Rolling variance
        pl.col("random") \
            .rolling_var(window_size=2) \
            .alias("rolling_var"),

        # Rolling quantile
        pl.col("random") \
            .rolling_quantile(
                quantile=0.75, window_size=2,
                min_periods=2
            ) \
            .alias("rolling_quantile"),

        # Rolling skew
        pl.col("random") \
            .rolling_skew(window_size=2) \
            .alias("rolling_skew"),

        # Rolling custom function
        pl.col("random") \
            .rolling_apply(
                function=np.nanstd, window_size=2) \
            .alias("rolling_apply")])
```

---

□□□□□□□□

□□□□□□□□□□□□□□□ □□□□□□□□□□□□□□□□□□□□

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

□□□□□□□□□□□



□□□□

□□□□□□□□□□□□□□□□□□□□

```python
df.join(df4, on="nrs")
# or
df.join(df4, on="nrs", how="inner")
```



□□□□

"left" □□□□ (□□□□□□□□) □□□□□□□□□□□□

```python
df.join(df4, on="nrs", how="left")
```



□□□□

□□□□□□□□□□□□□□□□□□□□□□□

```python
df.join(df4, on="nrs", how="outer")
```



□□□□□

df □□□ df4 □□□□□□□□□□□□□□□□□

```python
df.join(df4, on="nrs", how="anti")
```