# PySpark Functions Cheatsheet

## Column & Row Operations

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| withColumn() | Add/modify a column | df.withColumn("new", df["col"] + 1) |
| select() | Select specific columns | df.select("col1", "col2") |
| filter() / where() | Filter rows | df.filter(df["age"] > 30) |
| drop() | Drop a column | df.drop("col") |
| dropDuplicates() | Remove duplicates | df.dropDuplicates(["col"]) |
| distinct() | Unique rows | df.distinct() |
| withColumnRenamed() | Rename column | df.withColumnRenamed("old", "new") |
| orderBy() / sort() | Sort rows | df.orderBy("col", ascending=False) |
| cache() / persist() | Store DataFrame in memory/disk | df.cache() |

## Join & Combine

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| join() | Join two DataFrames | df1.join(df2, "id", "inner") |
| union() | Append DataFrames | df1.union(df2) |
| unionByName() | Append by matching column names | df1.unionByName(df2) |
| repartition() | Increase partitions | df.repartition(10) |
| coalesce() | Reduce partitions | df.coalesce(1) |

# Aggregation & Grouping

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| GROUPBY() | Group by column(s) | DF.GROUPBY("COL") |
| AGG() | Multiple aggregations | DF.AGG(F.MAX("COL"), F.MIN("COL")) |
| DESCRIBE() | Summary stats | DF.DESCRIBE().SHOW() |
| COLLECT_LIST() | List aggregation | COLLECT_LIST("NAME") |
| COLLECT_SET() | Set (distinct) aggregation | COLLECT_SET("NAME") |

# Null & Conditional Handling

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| FILLNA() | Replace nulls | DF.FILLNA(0) |
| ISNULL() | Check nulls | DF.FILTER(DF["COL"].ISNULL()) |
| WHEN() | If-else logic | WHEN(DF.AGE < 18, "MINOR").OTHERWISE("ADULT") |
| LIT() | Add literal value | LIT("USA") |

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| ROW_NUMBER() | Row number per window | ROW_NUMBER().OVER(WINDOW.PARTITIONBY("COL")) |
| RANK() | Rank with gaps | RANK().OVER(WINDOW.PARTITIONBY("COL")) |

| | | |
|---|---|---|
| DENSE_RANK() | RANK WITHOUT GAPS | DENSE_RANK().OVER(WINDOW.PARTITIONBY("COL")) |

# Window & Ranking Functions

# Date & Time

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| CURRENT_DATE() | TODAY'S DATE | CURRENT_DATE() |
| CURRENT_TIMESTAMP() | CURRENT TIMESTAMP | CURRENT_TIMESTAMP() |
| TO_DATE() | CONVERT TO DATE | TO_DATE(DF["TS"]) |
| DATE_FORMAT() | FORMAT DATE | DATE_FORMAT("DATE", "YYYY-MM") |

# Complex Types & JSON

| FUNCTION | DESCRIPTION | EXAMPLE |
|---|---|---|
| EXPLODE() | FLATTEN ARRAYS | EXPLODE("SKILLS") |
| ARRAY() / STRUCT() | CREATE COMPLEX TYPES | ARRAY("C1", "C2"), STRUCT("C1", "C2") |
| GET_JSON_OBJECT() | EXTRACT JSON FIELD | GET_JSON_OBJECT("JSON_COL", "$.NAME") |
| FROM_JSON() / TO_JSON() | PARSE OR STRINGIFY JSON | FROM_JSON("COL", SCHEMA) |