



PySpark Cheatsheet 2.0

[Explore More](#)

Contents

- Setup
- I/O
- Schema
- Cleaning
- Filtering
- Sorting
- Columns
- Strings
- Joins
- Aggregations
- Dates
- Window functions
- Arrays and Structs
- UDFs
- SCD
- Delta
- Streaming
- Performance
- Debugging

Part 1: Connect & Read

Import



```
1 # 1. Main Spark Session and Types
2 from pyspark.sql import SparkSession
3 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, DoubleType, DateType, TimestampType
4
5 # 2. Functions (Standard Alias 'F' is industry best practice)
6 from pyspark.sql import functions as F
7
8 # 3. Window Functions (for ranking, running totals)
9 from pyspark.sql.window import Window
10
11 # 4. Initialize Spark Session
12 spark = SparkSession.builder \
13     .appName("MasterCheatSheet") \
14     .config("spark.sql.adaptive.enabled", "true") \
15     .getOrCreate()
16
17 # 5. Check Spark Version
18 print(spark.version)
19
20 # 6. Set Log Level (Reduce noise in console)
21 spark.sparkContext.setLogLevel("WARN")
```

I/O

Schema



```
1 # 1. Define a Schema explicitly
2 schema = StructType([
3     StructField("id", IntegerType(), False),      # Not Nullable
4     StructField("name", StringType(), True),       # Nullable
5     StructField("salary", DoubleType(), True),
6     StructField("hire_date", DateType(), True)
7 ])
8
9 # 2. Apply Schema while reading (Fastest Read Method)
10 df = spark.read.schema(schema).csv("employees.csv")
11
12 # 3. Print Schema to console (Tree format)
13 df.printSchema()
14
15 # 4. Get List of Column Names
16 columns = df.columns
17
18 # 5. Inspect Column Data Types
19 dtypes = df.dtypes # Returns list of tuples: [('id', 'int'), ('name', 'string')]
20
21 # 6. Alter Schema: Change Nullability (Advanced)
22 # Creates a new schema based on existing one but forces fields to be Nullable
23 new_schema = StructType([
24     StructField(f.name, f.dataType, True) for f in df.schema.fields
25 ])
```

Part 2: Clean & Transform

Nulls & Duplicates



```
1 # 1. Drop Duplicate Rows (Checks all columns)
2 df = df.dropDuplicates()
3
4 # 2. Drop Duplicates based on specific columns (e.g., only check ID)
5 df = df.dropDuplicates(["id"])
6
7 # 3. Drop Rows containing ANY null values
8 df = df.na.drop(how="any")
9
10 # 4. Drop Rows where ALL columns are null
11 df = df.na.drop(how="all")
12
13 # 5. Fill Nulls with a default value (e.g., 0 for integers)
14 df = df.na.fill(0)
15
16 # 6. Fill Nulls for specific columns using a dictionary
17 df = df.na.fill({"age": 0, "city": "Unknown", "salary": 1000.0})
18
19 # 7. Replace specific values (e.g., empty strings with None)
20 df = df.replace("", None)
21
22 # 8. Trim Whitespace from strings (Left and Right)
23 df = df.withColumn("name", F.trim(F.col("name")))
24
25 # 9. Standardize Case (Lower/Upper)
26 df = df.withColumn("email", F.lower(F.col("email")))
```

Filtering



```
1 # 1. Simple Equality Filter
2 df_filt = df.filter(F.col("status") == "Active")
3
4 # 2. Multiple Conditions (AND) - Note: Parentheses are mandatory in PySpark
5 df_filt = df.filter((F.col("age") > 25) & (F.col("dept") == "Sales"))
6
7 # 3. Multiple Conditions (OR)
8 df_filt = df.filter((F.col("category") == "A") | (F.col("category") == "B"))
9
10 # 4. Filter using a List (isin)
11 df_filt = df.filter(F.col("country").isin("USA", "UK", "Canada"))
12
13 # 5. Filter for Null / Not Null
14 df_filt = df.filter(F.col("manager_id").isNotNull())
15 df_nulls = df.filter(F.col("manager_id").isNull())
16
17 # 6. Filter Strings (Contains / StartsWith)
18 df_filt = df.filter(F.col("email").contains("@gmail.com"))
19 df_filt = df.filter(F.col("name").startswith("J"))
20
21 # 7. Filter using Negation (~)
22 df_filt = df.filter(~(F.col("status") == "Closed"))
23
24 # 8. Filter Between Range
25 df_filt = df.filter(F.col("salary").between(50000, 100000))
```

Ordering

Column Operations



```
1 # 1. Add a New Column (Constant Value)
2 df = df.withColumn("report_date", F.lit("2023-01-01"))
3
4 # 2. Add a Calculated Column
5 df = df.withColumn("total_comp", F.col("salary") + F.col("bonus"))
6
7 # 3. Rename a Column
8 df = df.withColumnRenamed("fname", "first_name")
9
10 # 4. Conditional Logic (Case When ... Otherwise)
11 df = df.withColumn("grade",
12     F.when(F.col("score") >= 90, "A")
13     .when(F.col("score") >= 80, "B")
14     .otherwise("C")
15 )
16
17 # 5. Cast Data Type (Change String to Integer, etc.)
18 df = df.withColumn("age", F.col("age").cast(IntegerType()))
19
20 # 6. Concatenate Columns (Combine First + Last Name)
21 df = df.withColumn("full_name", F.concat_ws(" ", F.col("first"), F.col("last")))
22
23 # 7. Drop a Column
24 df = df.drop("temp_column")
25
26 # 8. Select specific columns to keep (Reorder)
27 df = df.select(F.col("id"), F.col("full_name"), F.col("grade"))
```

String Operations



```
1 # 1. Change Case (Lower / Upper / Title)
2 df = df.withColumn("lower_name", F.lower(F.col("name")))
3 df = df.withColumn("init_cap", F.initcap(F.col("description")))
4
5 # 2. Trim Whitespace (Both ends, Left, Right)
6 df = df.withColumn("clean_id", F.trim(F.col("id_str")))
7
8 # 3. Concatenate Strings (with Separator)
9 # Combines columns with a separator (ignores nulls automatically)
10 df = df.withColumn("full_address", F.concat_ws(", ", F.col("city"), F.col("state")))
11
12 # 4. Substring (Extract parts of a string)
13 # Note: PySpark substring index starts at 1, and length is the second arg
14 df = df.withColumn("zip_prefix", F.substring(F.col("zipcode"), 1, 3))
15
16 # 5. Split String (Convert String to Array)
17 df = df.withColumn("tags_array", F.split(F.col("tags"), ","))
18
19 # 6. Regex Replace (Pattern Matching)
20 # Remove all non-digit characters (keep only numbers)
21 df = df.withColumn("phone_clean", F regexp_replace(F.col("phone"), "[^0-9]", ""))
22
23 # 7. Check String Length
24 df = df.withColumn("desc_len", F.length(F.col("description")))
25
26 # 8. Check if String Contains Pattern
27 df = df.withColumn("is_gmail", F.col("email").contains("gmail"))
28
29 # 9. Padding (Add characters to reach fixed width)
30 df = df.withColumn("id_padded", F.lpad(F.col("id"), 10, "0"))
```

Part 3: Joins & Aggregate

Joins

Aggregations



```
1 # 1. Simple Count (Rows per Group)
2 df_agg = df.groupBy("department").count()
3
4 # 2. Multiple Aggregations (Sum, Avg, Max, Min)
5 # Always use .agg() to perform multiple calcs at once
6 df_agg = df.groupBy("department").agg(
7     F.sum("salary").alias("total_salary"),
8     F.avg("salary").alias("avg_salary"),
9     F.max("bonus").alias("max_bonus"),
10    F.count("id").alias("emp_count")
11 )
12
13 # 3. Count Distinct (Unique values per group)
14 df_agg = df.groupBy("category").agg(F.countDistinct("product_id"))
15
16 # 4. Pivot (Convert Rows to Columns)
17 # Example: Turn 'month' values (Jan, Feb) into columns
18 df_pivot = df.groupBy("product").pivot("month").sum("sales")
19
20 # 5. Global Aggregation (No GroupBy - summarize entire DF)
21 max_val = df.agg(F.max("salary")).collect()[0][0]
22
23 # 6. Collect List (Aggregate values into an Array)
24 # Warning: Be careful with large lists, can cause OOM errors
25 df_agg = df.groupBy("dept").agg(F.collect_list("name").alias("employees"))
```

Date & Time



```
1 # 1. Get Current Date & Timestamp
2 df = df.withColumn("now", F.current_timestamp())
3 df = df.withColumn("today", F.current_date())
4
5 # 2. Convert String to Date (Cast or Parse)
6 df = df.withColumn("date_col", F.to_date(F.col("date_str"), "yyyy-MM-dd"))
7
8 # 3. Convert Timestamp to String (Format)
9 df = df.withColumn("date_str", F.date_format(F.col("timestamp"), "MM/dd/yyyy"))
10
11 # 4. Date Arithmetic (Add / Subtract Days)
12 df = df.withColumn("next_week", F.date_add(F.col("date_col"), 7))
13 df = df.withColumn("last_week", F.date_sub(F.col("date_col"), 7))
14
15 # 5. Date Difference (Returns number of days)
16 df = df.withColumn("days_diff", F.datediff(F.col("end_date"), F.col("start_date")))
17
18 # 6. Extract Date Components
19 df = df.withColumn("year", F.year(F.col("date_col")))
20 df = df.withColumn("month", F.month(F.col("date_col")))
21 df = df.withColumn("day", F.dayofmonth(F.col("date_col")))
22
23 # 7. Truncate Date (Start of Month/Year)
24 # Useful for monthly reporting aggregations
25 df = df.withColumn("start_of_month", F.trunc(F.col("date_col"), "Month"))
26
27 # 8. Unix Timestamp (Seconds since Epoch)
28 df = df.withColumn("unix_time", F.unix_timestamp(F.col("timestamp_col")))
```

Window Functions



```
1 from pyspark.sql.window import Window
2
3 # 1. Define Window Specification (Partition & Order)
4 w_spec = Window.partitionBy("department").orderBy(F.col("salary").desc())
5
6 # 2. Row Number (Sequential: 1, 2, 3 ... )
7 df = df.withColumn("row_num", F.row_number().over(w_spec))
8
9 # 3. Rank (Gaps for ties: 1, 2, 2, 4)
10 df = df.withColumn("rank_val", F.rank().over(w_spec))
11
12 # 4. Dense Rank (No gaps for ties: 1, 2, 2, 3)
13 df = df.withColumn("dense_rank", F.dense_rank().over(w_spec))
14
15 # 5. Lead (Next Row's Value) - Shift Up
16 df = df.withColumn("next_salary", F.lead("salary", 1).over(w_spec))
17
18 # 6. Lag (Previous Row's Value) - Shift Down
19 df = df.withColumn("prev_salary", F.lag("salary", 1).over(w_spec))
20
21 # 7. Running Total (Cumulative Sum)
22 df = df.withColumn("running_sum", F.sum("salary").over(w_spec))
23
24 # 8. Group Average (Append average to every row - No OrderBy needed)
25 w_agg = Window.partitionBy("department")
26 df = df.withColumn("dept_avg", F.avg("salary").over(w_agg))
```

Part 4: Master & Extend

Complex Data Types



```
1 # 1. Explode Array (Create new row for each element)
2 # [A, B] → Row 1: A, Row 2: B
3 df = df.withColumn("tag", F.explode(F.col("tags_array")))
4
5 # 2. Array Contains (Filter if element exists)
6 df_filt = df.filter(F.array_contains(F.col("tags_array"), "urgent"))
7
8 # 3. Get Array Size
9 df = df.withColumn("num_tags", F.size(F.col("tags_array")))
10
11 # 4. Get Specific Element (1-based index)
12 df = df.withColumn("first_tag", F.element_at(F.col("tags_array"), 1))
13
14 # 5. Create Map (Key-Value pairs from columns)
15 df = df.withColumn("properties", F.create_map(F.lit("color"), F.col("color_col")))
16
17 # 6. Get Map Keys and Values
18 df = df.withColumn("all_keys", F.map_keys(F.col("properties")))
19 df = df.withColumn("all_values", F.map_values(F.col("properties")))
20
21 # 7. Access Nested Struct Field (Dot Notation)
22 # Assuming 'address' is a struct column with field 'city'
23 df = df.withColumn("city", F.col("address.city"))
24
25 # 8. Collapse Columns into Struct
26 df = df.withColumn("full_address", F.struct(F.col("city"), F.col("state")))
27
28 # 9. Collect Set (Aggregation: Create Array of Unique Values)
29 df_agg = df.groupBy("dept").agg(F.collect_set("id").alias("id_list"))
```

UDFs



```
1 from pyspark.sql.types import StringType, DoubleType
2 from pyspark.sql.functions import udf, pandas_udf
3 import pandas as pd
4
5 # 1. Standard Python UDF (Slow - Avoid for large data)
6 def upper_case_fn(s):
7     return s.upper() if s else None
8
9 # Register as UDF with return type
10 upper_udf = udf(upper_case_fn, StringType())
11 df = df.withColumn("upper_name", upper_udf(F.col("name")))
12
13 # 2. Pandas UDF (Vectorized - Much Faster)
14 # Input: pd.Series, Output: pd.Series
15 @pandas_udf(DoubleType())
16 def add_one_udf(a: pd.Series) → pd.Series:
17     return a + 1
18
19 df = df.withColumn("score_plus_one", add_one_udf(F.col("score")))
20
21 # 3. Register UDF for SQL Usage
22 spark.udf.register("sql_upper", upper_case_fn, StringType())
23 df_sql = spark.sql("SELECT sql_upper(name) FROM employees")
24
25 # 4. Apply UDF with Multiple Columns
26 @pandas_udf(StringType())
27 def combine_names(first: pd.Series, last: pd.Series) → pd.Series:
28     return first + " " + last
29
30 df = df.withColumn("full_name", combine_names(F.col("first"), F.col("last")))
```

SCD Type 2



```
1 # 1. Filter for active records in the target table
2 active = target.filter(F.col("is_active") == "Y")
3
4 # 2. Rename source columns to avoid ambiguity in join
5 src = source.select(*(F.col(c).alias(f"src_{c}") for c in source.columns))
6
7 # 3. Join Source with Active Target (Full Outer to catch all cases)
8 joined = active.join(src, F.col("id") == F.col("src_id"), "full")
9
10 # 4. Identify Updates (Match found BUT values changed)
11 # Note: Using strict inequality ( $\neq$ ) to detect changes
12 updates = joined.filter(
13     F.col("id").isNotNull() & F.col("src_id").isNotNull() &
14     (F.col("salary") < > F.col("src_salary")))
15 )
16
17 # 5. Identify New Records (No match in target)
18 new_recs = joined.filter(F.col("id").isNull())
19
20 # 6. Close Old Records (Expire: Set is_active=N, end_date=Today)
21 closed = updates.select("id", "salary", "start_date") \
22     .withColumn("end_date", F.current_date()) \
23     .withColumn("is_active", F.lit("N"))
24
25 # 7. Prepare New Active Records (From Updates + Brand New IDs)
26 # Combine, map back to original column names, and add active flags
27 active_new = updates.unionByName(new_recs) \
28     .select("src_id", "src_salary") \
29     .toDF("id", "salary") \
30     .withColumn("start_date", F.current_date()) \
31     .withColumn("end_date", F.lit(None)) \
32     .withColumn("is_active", F.lit("Y"))
33
34 # 8. Union All: Unchanged History + Closed Records + New Active
35 final_df = target.join(updates, "id", "left_anti") \
36     .unionByName(closed).unionByName(active_new)
```

Time Travel



```
1 from delta.tables import *
2
3 # 1. Read a specific version (e.g., Version 5)
4 df_ver = spark.read.format("delta").option("versionAsOf", 5).load("path/to/table")
5
6 # 2. Read by Timestamp (e.g., Data as it looked yesterday)
7 df_time = spark.read.format("delta") \
8     .option("timestampAsOf", "2023-10-25 12:00:00") \
9     .load("path/to/table")
10
11 # 3. Initialize DeltaTable Object (Required for maintenance ops)
12 delta_tbl = DeltaTable.forPath(spark, "path/to/table")
13
14 # 4. View History (See who changed what and when)
15 delta_tbl.history().show()
16
17 # 5. Restore to Version (Rollback - Overwrites current state!)
18 delta_tbl.restoreToVersion(5)
19
20 # 6. Restore to Timestamp
21 delta_tbl.restoreToTimestamp("2023-10-25")
22
23 # 7. SQL Time Travel (If table is registered in Catalog)
24 df_sql = spark.sql("SELECT * FROM employees TIMESTAMP AS OF '2023-10-25'")
25
26 # 8. Vacuum (Cleanup old files to save space)
27 # Warning: You cannot time travel back past the retention period (default 7 days)
28 # 168 hours = 7 days
29 delta_tbl.vacuum(168)
```

Streaming



```
1 # 1. Read Stream (Input)
2 # Requires a defined schema for file sources (JSON/CSV)
3 df_stream = spark.readStream.schema(schema).json("path/to/input_folder")
4
5 # 2. Transformations (Same as Batch!)
6 df_processed = df_stream.filter(F.col("status") == "error")
7
8 # 3. Write Stream to Console (Debugging only)
9 query = df_processed.writeStream \
10     .outputMode("append") \
11     .format("console") \
12     .start()
13
14 # 4. Write Stream to Delta (Production)
15 # Checkpointing is mandatory for fault tolerance
16 query = df_processed.writeStream \
17     .format("delta") \
18     .outputMode("append") \
19     .option("checkpointLocation", "path/to/checkpoints") \
20     .start("path/to/output_table")
21
22 # 5. Trigger (Control batch interval)
23 # "availableNow=True" processes all available data then stops (Batch-like)
24 query = df_processed.writeStream \
25     .trigger(processingTime="10 seconds") \
26     .start()
27
28 # 6. Watermarking (Handle late data in aggregations)
29 # Drops data older than the threshold (e.g., 10 minutes)
30 df_agg = df_stream \
31     .withWatermark("timestamp", "10 minutes") \
32     .groupBy(F.window("timestamp", "5 minutes")) \
33     .count()
```

Part 5: Optimize & Debug

Optimisation



```
1 from pyspark import StorageLevel
2
3 # 1. Cache (Lazy: Stores in memory after first action)
4 df.cache()
5
6 # 2. Persist (Memory & Disk - prevents OutOfMemory errors)
7 df.persist(StorageLevel.MEMORY_AND_DISK)
8
9 # 3. Unpersist (Clear memory when done - Critical!)
10 df.unpersist()
11
12 # 4. Repartition (Increases partitions, Full Shuffle)
13 # Use when you need more parallelism or to fix skewed data
14 df_repart = df.repartition(100)
15
16 # 5. Coalesce (Decreases partitions, No Shuffle)
17 # Use before writing to reduce number of output files
18 df_coal = df.coalesce(1)
19
20 # 6. Broadcast Join (For small tables < 10MB)
21 # Sends copy of small table to all nodes, avoids shuffle
22 df_join = df_large.join(F.broadcast(df_small), "id")
23
24 # 7. Explain Plan (Check physical execution)
25 df.explain()
26
27 # 8. Tune Shuffle Partitions (Default is 200)
28 # Set equal to 2x-3x total cores, or lower for small data
29 spark.conf.set("spark.sql.shuffle.partitions", "200")
```

Debugging



```
1 # 1. Explain Plan (Crucial for performance debugging)
2 # Shows the Physical Plan (indexes, join strategies used)
3 df.explain(mode="extended")
4
5 # 2. Show Full Content (Don't truncate long strings)
6 df.show(truncate=False, n=5)
7
8 # 3. Check for Skew (Count records per partition)
9 # If one partition has significantly more data, you have skew.
10 partition_counts = df.rdd.glom().map(len).collect()
11 print(partition_counts)
12
13 # 4. Identify Data Skew on a Key
14 # If one key has millions of records, joins will be slow.
15 df.groupBy("join_key").count().orderBy(F.col("count").desc()).show()
16
17 # 5. Check Data Types (List format)
18 print(df.dtypes)
19
20 # 6. Validate Nulls in Critical Columns
21 null_counts = df.select([
22     F.count(F.when(F.col(c).isNull(), c)).alias(c)
23     for c in df.columns
24 ])
25 null_counts.show()
26
27 # 7. Clear Cache (If you run out of memory during iterative dev)
28 spark.catalog.clearCache()
29
30 # 8. Check Spark Configuration (Debug memory settings)
31 print(spark.conf.get("spark.sql.shuffle.partitions"))
```