

Spark Cheat Sheet

Start Spark & Load Data

Load CSV

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("MyApp").getOrCreate()
df = spark.read.csv("filename.csv", header=True, inferSchema=True)
```

Load JSON

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("MyApp").getOrCreate()
df = spark.read.json("filename.json")
```

Exploring the Data

Task	Command
Show rows	<code>df.show()</code>
Show first N rows	<code>df.show(10)</code>
Show schema	<code>df.printSchema()</code>
Show column names	<code>df.columns</code>
Get number of rows	<code>df.count()</code>
Get distinct values	<code>df.select("column").distinct().show()</code>
Describe (summary stats)	<code>df.describe().show()</code>

Column operations

Task	Command
Select column(s)	<code>df.select("col1", "col2")</code>
Filter rows	<code>df.filter(df.col1 > 50)</code>
Add new column	<code>df.withColumn("new_col", col("col1") + 5)</code>
Rename column	<code>df.withColumnRenamed("old", "new")</code>
Cast column type	<code>df.withColumn("col", df["col"].cast("float"))</code>
Drop a column	<code>df.drop("col_to_drop")</code>

Row operations

Task	Command
Show first N rows	<code>df.show(N)</code>
Show top rows with full content	<code>df.show(truncate=False)</code>
Filter rows with a condition	<code>df.filter(df["col"] > 100)</code>
Filter with multiple conditions	<code>df.filter((df["age"] > 18) & (df["country"] == "US"))</code>
Filter using <code>where()</code>	<code>df.where("age > 30")</code>
Filter NULL rows	<code>df.filter(df["col"].isNull())</code>
Filter non-NULL rows	<code>df.filter(df["col"].isNotNull())</code>
Select rows with values in list	<code>df.filter(df["col"].isin("A", "B", "C"))</code>
Drop duplicate rows	<code>df.dropDuplicates()</code>
Drop duplicate rows by column(s)	<code>df.dropDuplicates(["col1", "col2"])</code>
Limit number of rows	<code>df.limit(100)</code>
Random sample of rows	<code>df.sample(withReplacement=False, fraction=0.1)</code>
Add row number	<code>from pyspark.sql.window import Window`from pyspark.sql.functions import row_number`df.withColumn("row_num", row_number().over(Window.orderBy("col")))</code>
Add conditional column	<code>df.withColumn("flag", when(df["score"] > 70, "Pass").otherwise("Fail"))</code>
Replace values conditionally	<code>df.withColumn("status", when(df["status"] == "old", "new").otherwise(df["status"]))</code>

Data cleaning

Task	Command
Drop rows with nulls	<code>df.dropna()</code>
Drop nulls in specific columns	<code>df.dropna(subset=["col1", "col2"])</code>
Fill nulls	<code>df.fillna({'col1': 0, 'col2': 'Unknown'})</code>
Replace values	<code>df.replace("old", "new", "column")</code>
Drop duplicates	<code>df.dropDuplicates()</code>
Drop duplicates by column	<code>df.dropDuplicates(["col1"])</code>
Trim strings	<code>from pyspark.sql.functions import trim; df.withColumn("col", trim(col("col")))</code>
Convert to lower case	<code>df.withColumn("col", lower(col("col")))</code>
Convert to upper case	<code>df.withColumn("col", upper(col("col")))</code>
Remove special characters	<code>df.withColumn("col", regexp_replace(col("col"), "[^a-zA-Z0-9]", ""))</code>
Cast column type	<code>df.withColumn("col", col("col").cast("float"))</code>
Rename column	<code>df.withColumnRenamed("old_name", "new_name")</code>
Split string into array	<code>df.withColumn("col_array", split(col("col"), ","))</code>
Explode array into rows	<code>df.withColumn("value", explode(col("col_array")))</code>
Add column with condition	<code>df.withColumn("flag", when(col("score") > 90, "high").otherwise("low"))</code>
Parse date	<code>df.withColumn("parsed_date", to_date("Date", "M/d/yyyy"))</code>
Extract year/month/day from date	<code>df.withColumn("Year", year(col("parsed_date"))) etc.</code>
Check for NULLs in all columns	<code>df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns]).show()</code>
Select rows with NULLs	<code>df.filter(col("col").isNull())</code>
Select rows with non-NULLs	<code>df.filter(col("col").isNotNull())</code>
Conditional replacement	<code>df.withColumn("col", when(col("col") == "old", "new").otherwise(col("col")))</code>
Round numeric values	<code>df.withColumn("col_rounded", round(col("col"), 2))</code>

Aggregations

Task	Command
Group and count	<code>df.groupBy("col").count()</code>
Group and average	<code>df.groupBy("col").avg("value_col")</code>
Multiple aggregations	<code>df.groupBy("col").agg({'value_col': 'avg', 'id': 'count'})</code>
Sort ascending	<code>df.orderBy("col")</code>
Sort descending	<code>df.orderBy(df.col.desc())</code>

Using aggregated functions

```
from pyspark.sql.functions import avg, max, min, sum, stddev
df.select(avg("Valuation_numeric"), stddev("Valuation_numeric")).show()
```

Date functions

```
# For all options
from pyspark.sql.functions import (
    to_date, to_timestamp, year, month, dayofmonth, dayofweek, weekofyear,
    date_format, date_add, date_sub, datediff, concat_ws, current_date
)

from pyspark.sql.functions import to_date, year, month, dayofmonth

df = df.withColumn("DateParsed", to_date("Date", "M/d/yyyy"))
df = df.withColumn("Year", year("DateParsed"))
df = df.withColumn("Month", month("DateParsed"))
df = df.withColumn("Day", dayofmonth("DateParsed"))
```

Date/time operations

Task	Command
Convert string to date (MM/dd/yyyy)	<code>df.withColumn("parsed_date", to_date("Date", "MM/dd/yyyy"))</code>
Convert string to date (M/d/yyyy)	<code>df.withColumn("parsed_date", to_date("Date", "M/d/yyyy"))</code>
Convert to timestamp	<code>df.withColumn("ts", to_timestamp("Date", "MM/dd/yyyy HH:mm:ss"))</code>
Extract year	<code>df.withColumn("Year", year("parsed_date"))</code>
Extract month	<code>df.withColumn("Month", month("parsed_date"))</code>
Extract day of month	<code>df.withColumn("Day", dayofmonth("parsed_date"))</code>
Extract day of week (1=Sun)	<code>df.withColumn("DayOfWeek", dayofweek("parsed_date"))</code>
Extract week of year	<code>df.withColumn("WeekOfYear", weekofyear("parsed_date"))</code>
Format date to string (yyyy-MM)	<code>df.withColumn("MonthFormatted", date_format("parsed_date", "yyyy-MM"))</code>
Add days	<code>df.withColumn("Plus7Days", date_add("parsed_date", 7))</code>
Subtract days	<code>df.withColumn("Minus7Days", date_sub("parsed_date", 7))</code>
Calculate difference between dates	<code>df.withColumn("diff_days", datediff(col("end_date"), col("start_date")))</code>
Create date from year, month, day	<code>df.withColumn("full_date", to_date(concat_ws("-", "year", "month", "day")))</code>
Get current date/time	<code>from pyspark.sql.functions import current_date, current_timestamp`current_date()</code>

Save as `csv`

Task	Command
Save as CSV	<code>df.write.csv("folder", header=True)</code>
Save as single CSV file	<code>df.coalesce(1).write.csv("folder", header=True)</code>
Save as Parquet	<code>df.write.parquet("path")</code>

Transformations vs Actions

Category	Transformations	Actions
Definition	Define <i>what to do</i> with the data	Actually <i>do it</i> and return results
Lazy?	✅ Yes (do nothing until an action is called)	❌ No (they trigger execution)
⚙️Returns	New DataFrame	Actual values or saved output
Triggers execution?	No	Yes
Examples	<code>select()</code> , <code>filter()</code> , <code>withColumn()</code> , <code>groupBy()</code>	<code>show()</code> , <code>collect()</code> , <code>count()</code> , <code>take()</code> , <code>write()</code>

Transformations

Function	Description	Example Usage
<code>pivot()</code>	Rotates rows into columns (like a pivot table)	<code>df.groupBy("Year").pivot("Country").sum("Valuation")</code>
<code>explode()</code>	Converts an array or string split into multiple rows	<code>explode(split(col("Investors"), ","))</code>
<code>split()</code>	Splits a string into an array based on a delimiter	<code>split(col("Investors"), ",")</code>
<code>array()</code>	Combines multiple columns into an array	<code>array("col1", "col2")</code>
<code>concat_ws()</code>	Concatenates multiple columns into one string (with a separator)	<code>concat_ws("-", col("Year"), col("Month"), col("Day"))</code>
<code>withColumn()</code>	Creates or replaces a column	<code>df.withColumn("Val", col("Valuation").cast("float"))</code>
<code>selectExpr()</code>	SQL-style expression for selecting and transforming columns	<code>df.selectExpr("Company", "Valuation * 1.2 as AdjustedVal")</code>
<code>melt()</code> (custom)	Converts wide data back to long format (inverse of pivot — requires workaround)	[Not native, use stack/union workaround]
<code>groupBy().agg()</code>	Aggregate after grouping, with multiple functions	<code>df.groupBy("Industry").agg(avg("Val"), count("*"))</code>
<code>stack()</code>	Emulates melt/unpivot by converting multiple columns into rows	<code>SELECT Company, stack(2, '2020', val2020, '2021', val2021) AS (Year, Value) FROM table</code>
<code>flatten()</code>	Flattens nested arrays (requires Spark 2.4+)	<code>flatten(col("nested_array"))</code>
<code>explode_outer()</code>	Same as <code>explode()</code> , but returns NULL if input is NULL	<code>explode_outer(col("arr"))</code>
<code>arrays_zip()</code>	Combines multiple arrays into an array of structs (zip-like behavior)	<code>arrays_zip("arr1", "arr2")</code>
<code>posexplode()</code>	Explodes array with element positions (like index + value)	<code>df.select(posexplode(split(col("tags"), ",")))</code>
<code>map()</code> (in SQL)	Create a map column from key-value pairs	<code>select map("k1", col1, "k2", col2) as kv_map</code>

Actions

Action	Description**	Returns
<code>show()</code>	Displays rows from the DataFrame	None (prints to console)
<code>collect()</code>	Collects all rows to the driver node	Python list of Rows
<code>take(n)</code>	Returns the first <code>n</code> rows	List of Row objects
<code>first()</code>	Returns the first row	Single Row
<code>head()</code>	Same as <code>first()</code> or <code>take(1)</code>	Row or list
<code>count()</code>	Counts number of rows	Integer
<code>describe().show()</code>	Summary statistics of numerical columns	Console output
<code>agg().show()</code>	Performs aggregate computations (avg, sum, etc.)	None (prints result)
<code>foreach()</code>	Runs a function on each row (no return)	None
<code>foreachPartition()</code>	Runs a function on each partition	None
<code>toPandas()</code>	Converts DataFrame to Pandas DataFrame	<code>pandas.DataFrame</code>
<code>write.csv(...)</code>	Saves DataFrame to CSV	Saves files (returns nothing)
<code>write.json(...)</code>	Saves DataFrame to JSON	Saves files
<code>write.parquet(...)</code>	Saves DataFrame to Parquet	Saves files
<code>write.saveAsTable(...)</code>	Saves DataFrame as a Hive table (if using Hive)	None
<code>isEmpty()</code> (3.3+)	Checks if DataFrame has no rows	Boolean
<code>inputFiles()</code> (<i>rdd</i>)	Returns paths used to construct a DataFrame from files	List of strings

Pivot

Pivoting is like rotating a table: it turns **row values into column headers**. It's useful for summarizing and reshaping data (like pivot tables in Excel).

```
df.groupBy("Year").pivot("Country").sum("Valuation")
```

Explode

Explode transforms an **array or map column into multiple rows**. It's like flattening a list inside a column.

For example, going from:

Company	Investors
Stripe	["Sequoia", "Tiger Global"]
Klarna	["BlackRock", "Atomico"]

To:

Company	Investor
Stripe	Sequoia
Stripe	Tiger Global
Klarna	BlackRock
Klarna	Atomico

```
from pyspark.sql.functions import explode, split

df.withColumn("Investor", explode(split(col("Investors"), ",")))
```