

Collect()

• `collect()` is an Action that brings all rows of a DataFrame/RDD **to the driver** as a python list

• Executors → Drivers (everything at once)

• Spark runs distributed, but `collect()` pulls all data into one machine's memory

Syntax:

`rows = df.collect()`

→ returns: `List[row]`

→ triggers execution immediately [ACTION]

Ex:

```
rows = df.select('id', 'name').collect()  
rows[0]['name']
```

Why `collect()` is dangerous?

- OOM risk (driver runs out of memory)
- slow for large datasets
- crashes notebooks/jobs in Databricks
- Never use on large tables

When is `collect()` acceptable?

- use only when you're sure the data is tiny:
 - > Debugging a few rows
 - > unit tests
 - > fetching a small lookup list

Alternatives of `collect()`: `show(20)`, `take(1)`, `head(1)`, `limit(1)`

StructType & StructField

24.10.2023

StructType : defines the entire schema (table structure)

StructField : defines one column inside that schema

- Together, they let us explicitly control column names, data types & nullability.

why we should code?

- If we don't define a schema, then Spark infers it
- Inference can be slow, wrong, inconsistent
- Explicit schema → faster + safer + production-ready

```
from pyspark.sql.types import (
```

```
    StructType, StructField,
```

```
    IntegerType, StringType, DoubleType
```

```
)
```

StructField("id", IntegerType, nullable=True)

```
schema = StructType([
```

```
    StructField("id", IntegerType, True),
```

```
    StructField("Name", StringType, True)
```

```
)
```

→ creating Dataframe with explicit schema

```
data = [
```

```
(1, "A"), (2, "B")]
```

`df = spark.createDataFrame(data, schema)`

* Reading files with schema:

`df = spark.read \`

- `schema(schema)`
- `option('header', True)`
- `csv(—)`

- Faster than inference
- Prevents wrong types

Nested schema:

`nestedSchema = StructType([`

`StructField("id", IntegerType(), True),`

`StructField("address", StructType([`

`StructField("city", StringType(), True),`

`StructField("zip", StringType(), True)`

`]), True)`

`])`

{
 ⇒ `RDD = Dataframe with StructType`

Pivot() & Unpivot()

pivot() → turns row values into columns

unpivot() → turns columns back into rows

Pivot() { row → columns }

When to use :

- > reports / dashboards
- > category wise metrics as columns
- > SQL-style crosstabs

Basic Syntax :



df.groupby(group-cols).pivot(pivot-col).agg(agg-fn)

Ex:

State	Product	Sales
AP	Shoes	100
AP	Shirts	80
TS	Shoes	120

pivot product-wise sales per state:

```
df.groupby("state")\n    .pivot("product")\n    .agg(sum('sales'))\n
```

Op

State	Shoes	Shirts
AP	100	80
TS	120	NULL

→ If we know pivot values, specify them to avoid an extra pass:

↓

```
df.groupBy('state')  
    .pivot("product", ["shoes", "shirts"])  
    .agg(sum("sales"))
```

→ Common rules:

- ↓
- ① `Pivot()` must be used with `groupBy()`
 - ② Always needs an aggregation
 - ③ Too many pivot values \Rightarrow Many columns
 \uparrow
(memory risk)

Unpivot(): { Columns \rightarrow Rows }

→ pyspark doesn't have a single `unpivot()` function everywhere,

but you can do it ~~cleanly~~ cleanly using `stack()`
(or SQL)

⇒ When to use:

- > Normalize wide data
- > ML features
- > Schema standardization

`df_unpivot = df.select(
 "state",
 expr("stack(2, 'shoe',
 'shoes', 'Shirts', Shirts) as
 (product, sales)"))`

→ we have;

state	shoes	shirts
AP	100	80
TS	120	NULL

we want →

state	product	sales
-------	---------	-------

That is; { column names become values }
{ column values become row values }

How unpivot works?

for each row;

take column "Shoes" → create a row (product:shoes,
sales = value)

take column "Shirts" → create another row

df.unpivot = df.select(

"state",

exprs(

After "State",

how many columns
we need to unpivot

"stack(2, "Shoes", "Shoes",
"Shirts", "Shirts")

as (product, sales) ?

- 2 → number of columns to unpivot
 'Shoes' → new value for product
 Shoes → value goes into sales
 'Shirts' → next product
 Shirts → next sales value

⇒ Why unpivot is useful?

- 1. ML & analysis
- 2. Schema consistency
- 3. reverse of pivot