

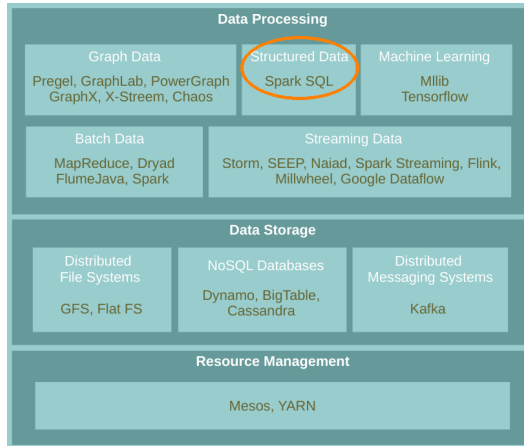


Structured Data Processing - Spark SQL

Amir H. Payberah
payberah@kth.se
2025-09-15



Where Are We?



Motivation

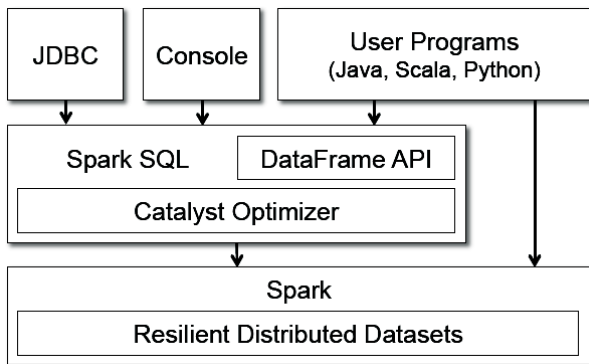
Structured Data



Unstructured Data



Spark and Spark SQL



Structured Data vs. RDD (1/2)

► `case class Account(name: String, balance: Double, risk: Boolean)`



Structured Data vs. RDD (1/2)

- ▶ `case class Account(name: String, balance: Double, risk: Boolean)`
- ▶ `RDD[Account]`



Structured Data vs. RDD (1/2)

- ▶ `case class Account(name: String, balance: Double, risk: Boolean)`
- ▶ `RDD[Account]`
- ▶ **RDDs don't know** anything about the **schema** of the data it's dealing with.





Structured Data vs. RDD (2/2)

- ▶ `case class Account(name: String, balance: Double, risk: Boolean)`
- ▶ `RDD[Account]`
- ▶ A **database/Hive** sees it as a columns of named and typed values.

name: String	balance: Double	risk: Boolean
name: String	balance: Double	risk: Boolean
name: String	balance: Double	risk: Boolean
name: String	balance: Double	risk: Boolean



DataFrames and DataSets

- ▶ Spark has **two** notions of **structured collections**:
 - **DataFrames**
 - **Datasets**
- ▶ They are **distributed table-like collections** with **well-defined rows and columns**.



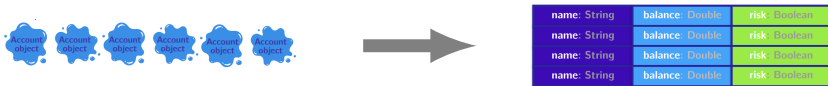
DataFrames and DataSets

- ▶ Spark has **two** notions of **structured collections**:
 - **DataFrames**
 - **Datasets**
- ▶ They are **distributed table-like collections** with **well-defined rows and columns**.
- ▶ They represent **immutable lazily** evaluated plans.
- ▶ When an **action** is performed on them, Spark performs the **actual transformations** and return the result.

DataFrame

DataFrame

- ▶ Consists of a **series of rows** and a **number of columns**.
- ▶ Equivalent to a **table** in a relational database.
- ▶ **Spark + RDD**: **functional** transformations on partitioned collections of **objects**.
- ▶ **SQL + DataFrame**: **declarative** transformations on partitioned collections of **tuples**.



- ▶ Defines the **column names and types** of a DataFrame.
- ▶ Assume **people.json** file as an input:

```
{ "name": "Michael", "age": 15, "id": 12 }  
{ "name": "Andy", "age": 30, "id": 15 }  
{ "name": "Justin", "age": 19, "id": 20 }  
{ "name": "Andy", "age": 12, "id": 15 }  
{ "name": "Jim", "age": 19, "id": 20 }  
{ "name": "Andy", "age": 12, "id": 10 }
```

- ▶ Defines the **column names and types** of a DataFrame.
- ▶ Assume `people.json` file as an input:

```
{ "name": "Michael", "age": 15, "id": 12 }  
{ "name": "Andy", "age": 30, "id": 15 }  
{ "name": "Justin", "age": 19, "id": 20 }  
{ "name": "Andy", "age": 12, "id": 15 }  
{ "name": "Jim", "age": 19, "id": 20 }  
{ "name": "Andy", "age": 12, "id": 10 }
```

```
val people = spark.read.format("json").load("people.json")  
people.schema
```

```
// returns:  
StructType(StructField(age,LongType,true),  
StructField(id,LongType,true),  
StructField(name,StringType,true))
```



Column

- ▶ They are like **columns in a table**.
- ▶ `col` returns a reference to a column.
- ▶ `columns` returns all columns on a DataFrame

```
val people = spark.read.format("json").load("people.json")  
  
col("age")  
  
people.columns  
// returns: Array[String] = Array(age, id, name)
```



Row

- ▶ A **row** is a **record of data**.
- ▶ They are of type **Row**.

```
import org.apache.spark.sql.Row  
  
val myRow = Row("Seif", 65, 0)
```




Row

- ▶ A **row** is a **record of data**.
- ▶ They are of type **Row**.
- ▶ Rows do **not have schemas**.

```
import org.apache.spark.sql.Row  
  
val myRow = Row("Seif", 65, 0)
```



Row

- ▶ A **row** is a **record of data**.
- ▶ They are of type **Row**.
- ▶ Rows do **not have schemas**.
- ▶ To access data in rows, you need to specify the **position** that you would like.

```
import org.apache.spark.sql.Row
```

```
val myRow = Row("Seif", 65, 0)
```

```
myRow(0) // type Any
```

```
myRow(0).asInstanceOf[String] // String
```

```
myRow.getString(0) // String
```

```
myRow.getInt(1) // Int
```



Creating a DataFrame

- ▶ **Two** ways to create a DataFrame:
 1. From an **RDD**
 2. From **raw data sources**



Creating a DataFrame - From an RDD

- ▶ You can use `toDF` to convert an RDD to DataFrame.
- ▶ The schema automatically inferred.

```
val tupleRDD = sc.parallelize(Array(("seif", 65, 0), ("amir", 40, 1)))  
val tupleDF = tupleRDD.toDF("name", "age", "id")
```



Creating a DataFrame - From an RDD

- ▶ You can use `toDF` to convert an RDD to DataFrame.
- ▶ The schema automatically inferred.

```
val tupleRDD = sc.parallelize(Array(("seif", 65, 0), ("amir", 40, 1)))  
val tupleDF = tupleRDD.toDF("name", "age", "id")
```

- ▶ If RDD contains `case` class instances, Spark infers the attributes from it.

```
case class Person(name: String, age: Int, id: Int)  
val peopleRDD = sc.parallelize(Array(Person("seif", 65, 0), Person("amir", 40, 1)))  
val peopleDF = peopleRDD.toDF
```



Creating a DataFrame - From Data Source

- Data sources supported by Spark.

```
val peopleJson = spark.read.format("json").load("people.json")

val peopleCsv = spark.read.format("csv")
  .option("sep", ";")
  .option("inferSchema", "true")
  .option("header", "true")
  .load("people.csv")
```



Creating a DataFrame - From Data Source

- ▶ Data sources supported by Spark.
 - CSV, JSON, Parquet, ORC, JDBC/ODBC connections, Plain-text files

```
val peopleJson = spark.read.format("json").load("people.json")

val peopleCsv = spark.read.format("csv")
  .option("sep", ";")
  .option("inferSchema", "true")
  .option("header", "true")
  .load("people.csv")
```



Creating a DataFrame - From Data Source

- ▶ Data sources supported by Spark.
 - CSV, JSON, Parquet, ORC, JDBC/ODBC connections, Plain-text files
 - Cassandra, HBase, MongoDB, AWS Redshift, XML, etc.

```
val peopleJson = spark.read.format("json").load("people.json")  
  
val peopleCsv = spark.read.format("csv")  
  .option("sep", ";")  
  .option("inferSchema", "true")  
  .option("header", "true")  
  .load("people.csv")
```


-
- Remove columns or rows
 - Transform a row into a column or a column into a row
 - Add rows or columns
 - Sort data by values in rows



DataFrame Transformations (2/5)

- ▶ `select` and `selectExpr` allow to do the **DataFrame equivalent** of **SQL queries** on a table of data.

```
// select  
people.select("name", "age", "id").show(2)
```



DataFrame Transformations (2/5)

- ▶ `select` and `selectExpr` allow to do the **DataFrame equivalent** of **SQL queries** on a table of data.

```
// select  
people.select("name", "age", "id").show(2)
```

```
// selectExpr  
people.selectExpr("*", "(age < 20) as teenager").show()
```



DataFrame Transformations (3/5)

- ▶ `filter` and `where` both `filter` rows.
- ▶ `distinct` can be used to extract unique rows.

```
people.filter("age < 20").show()
people.where("age < 20").show()
people.select("name").distinct().show()
```

What is the output?

```
people.selectExpr("avg(age)", "count(distinct(name)) as distinct").show()
```

```
+---+---+-----+
|age| id|   name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|   Andy|
| 19| 20|   Andy|
+---+---+-----+
```

What is the output?

```
people.selectExpr("avg(age)", "count(distinct(name)) as distinct").show()
```

```
+---+---+-----+
|age| id|   name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|   Andy|
| 19| 20|   Andy|
+---+---+-----+
```

```
+-----+-----+
|avg(age)|distinct|
+-----+-----+
|  21.333|        2|
+-----+-----+
```



DataFrame Transformations (4/5)

- ▶ `withColumn` adds a new column to a DataFrame.

```
people.withColumn("teenager", expr("age < 20")).show()
```



DataFrame Transformations (4/5)

- ▶ `withColumn` adds a new column to a DataFrame.
- ▶ `withColumnRenamed` renames a column.

```
people.withColumn("teenager", expr("age < 20")).show()
```

```
people.withColumnRenamed("name", "username").columns
```




DataFrame Transformations (4/5)

- ▶ `withColumn` adds a new column to a DataFrame.
- ▶ `withColumnRenamed` renames a column.
- ▶ `drop` removes a column.

```
people.withColumn("teenager", expr("age < 20")).show()
```

```
people.withColumnRenamed("name", "username").columns
```

```
people.drop("name").columns
```

What is the output?

```
people.withColumn("teenager", expr("age < 20")).show()
```

```
+---+---+-----+
|age| id|  name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|  Andy|
| 19| 20| Justin|
+---+---+-----+
```

What is the output?

```
people.withColumn("teenager", expr("age < 20")).show()
```

```
+---+---+-----+
|age| id|   name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|   Andy|
| 19| 20|  Justin|
+---+---+-----+
```

Option 1

```
+---+---+-----+
|age| id|   name|teenager|
+---+---+-----+
| 15| 12|Michael|    true|
| 30| 15|   Andy|   false|
| 19| 20|  Justin|    true|
+---+---+-----+
```

Option 2

```
+---+---+-----+
|age| id|   name|teenager|
+---+---+-----+
| 15| 12|Michael|    true|
| 19| 20|  Justin|    true|
+---+---+-----+
```



DataFrame Transformations (5/5)

- ▶ You can use `udf` to define new **column-based functions**.

```
import org.apache.spark.sql.functions.{col, udf}

val df = spark.createDataFrame(Seq((0, "hello"), (1, "world"))).toDF("id", "text")

val upper: String => String = _.toUpperCase
val upperUDF = spark.udf.register("upper", upper)

df.withColumn("upper", upperUDF(col("text"))).show
```



DataFrame Actions

- ▶ Like RDDs, DataFrames also have their own set of actions.
- ▶ `collect`: returns an `array` that contains all of `rows` in this DataFrame.
- ▶ `count`: returns the `number of rows` in this DataFrame.
- ▶ `first` and `head`: returns the `first row` of the DataFrame.
- ▶ `show`: displays the `top 20 rows` of the DataFrame in a tabular form.
- ▶ `take`: returns the `first n rows` of the DataFrame.

Aggregation



Aggregation

- ▶ In an **aggregation** you specify
 - A **key or grouping**
 - An **aggregation function**
- ▶ The given function must produce **one** result for **each group**.



Grouping Types

- ▶ Summarizing a complete DataFrame
- ▶ Group by
- ▶ Windowing



Grouping Types

- ▶ Summarizing a complete DataFrame
- ▶ Group by
- ▶ Windowing



Summarizing a Complete DataFrame Functions

- ▶ `count` returns the total number of values.

```
people.select(count("age")).show()
```



Summarizing a Complete DataFrame Functions

- ▶ `count` returns the total number of values.
- ▶ `countDistinct` returns the number of unique groups.

```
people.select(count("age")).show()
```

```
people.select(countDistinct("name")).show()
```



Summarizing a Complete DataFrame Functions

- ▶ `count` returns the **total number of values**.
- ▶ `countDistinct` returns the **number of unique groups**.
- ▶ `first`, `last`, `min`, `max`, `sum`, `avg`

```
people.select(count("age")).show()
```

```
people.select(countDistinct("name")).show()
```

```
people.select(first("name"), last("age"), min("age"), max("age"), sum("age"), avg("age")).show()
```



Grouping Types

- ▶ Summarizing a complete DataFrame
- ▶ Group by
- ▶ Windowing



Group By (1/2)

- ▶ Perform aggregations on **groups** in the data.
- ▶ Typically on **categorical data**.



Group By (1/2)

- ▶ Perform aggregations on **groups** in the data.
- ▶ Typically on **categorical data**.
- ▶ We do this grouping in **two phases**:
 1. **Specify the column(s)** on which we would like to group.
 2. Specify the **aggregation(s)** using the **agg** operation.



Group By (2/2)

- ▶ Grouping with expressions

```
people.groupBy("name").agg(count("age").alias("ageagg")).show()
```




Group By (2/2)

- ▶ Grouping with **expressions**
- ▶ Specifying transformations as a **series of Maps**

```
people.groupBy("name").agg(count("age").alias("ageagg")).show()
```

```
people.groupBy("name").agg("age" -> "count", "age" -> "avg", "id" -> "max").show()
```

What is the output?

```
people.groupBy("name").agg("age" -> "count", "age" -> "avg", "id" -> "max").show()
```

```
+---+---+-----+
|age| id|  name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|  Andy|
| 19| 20|  Andy|
+---+---+-----+
```

What is the output?

```
people.groupBy("name").agg("age" -> "count", "age" -> "avg", "id" -> "max").show()
```

```
+---+---+-----+
|age| id|  name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|  Andy|
| 19| 20|  Andy|
+---+---+-----+
```

Option 1

```
+-----+-----+-----+-----+
|  name|count(age)|avg(age)|max(id)|
+-----+-----+-----+-----+
|Michael|          1|   15.0|    12|
|  Andy|          2|   24.5|    20|
+-----+-----+-----+-----+
```

Option 2

```
+-----+-----+-----+-----+
|  name|count(age)|avg(age)|max(id)|
+-----+-----+-----+-----+
|Michael|          1|   21.33|    20|
|  Andy|          2|   21.33|    20|
+-----+-----+-----+-----+
```

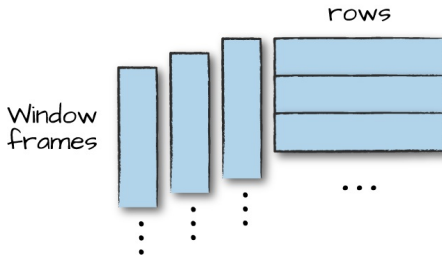


Grouping Types

- ▶ Summarizing a complete DataFrame
- ▶ Group by
- ▶ Windowing

Windowing (1/2)

- ▶ Computing some aggregation on a specific **window** of data.
- ▶ The **window** determines **which rows** will be passed in to this function.
- ▶ You define them by using a **reference to the current data**.



[M. Zaharia et al., Spark: The Definitive Guide, O'Reilly Media, 2018]

Windowing (2/2)

- Unlike grouping, here **each row** can fall into **one or more frames**.

```
import org.apache.spark.sql.expressions.Window
import org.apache.spark.sql.functions.col

val people = spark.read.format("json").load("people.json")

val windowSpec = Window.rowsBetween(-1, 1)
val avgAge = avg(col("age")).over(windowSpec)
people.select(col("name"), col("age"), avgAge.alias("avg_age")).show
```

What is the output?

```
val windowSpec = Window.rowsBetween(-1, 1)
val avgAge = avg(col("age")).over(windowSpec)
people.select(col("name"), col("age"), avgAge.alias("avg_age")).show()
```

```
+---+---+-----+
|age| id|  name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|  Andy|
| 19| 20|  Andy|
+---+---+-----+
```

What is the output?

```
val windowSpec = Window.rowsBetween(-1, 1)
val avgAge = avg(col("age")).over(windowSpec)
people.select(col("name"), col("age"), avgAge.alias("avg_age")).show()
```

```
+---+---+-----+
|age| id|   name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|   Andy|
| 19| 20|   Andy|
+---+---+-----+
```

Option 1

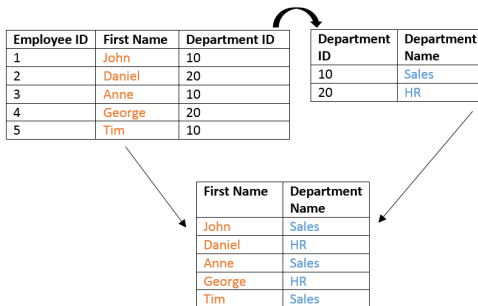
```
+-----+---+-----+
|   name|age| avg_age|
+-----+---+-----+
|Michael| 15|    22.5|
|   Andy| 30|    21.33|
|   Andy| 19|    24.5|
+-----+---+-----+
```

Option 2

```
+-----+---+-----+
|   name|age| avg_age|
+-----+---+-----+
|Michael| 15|     7.5|
|   Andy| 30|    22.5|
|   Andy| 19|    21.33|
+-----+---+-----+
```


Joins

- ▶ Joins are **relational** constructs you use to **combine relations together**.
- ▶ Different **join types**: inner join, outer join, left outer join, right outer join, left semi join, left anti join, cross join





Joins Example

```
val person = Seq((0, "Seif", 0), (1, "Amir", 1), (2, "Sarunas", 1))  
               .toDF("id", "name", "group_id")  
  
val group = Seq((0, "SICS/KTH"), (1, "KTH"), (2, "SICS"))  
              .toDF("id", "department")
```

Joins Example - Inner

```
val joinExpression = person.col("group_id") === group.col("id")  
  
var joinType = "inner"  
  
person.join(group, joinExpression, joinType).show()
```

```
+---+-----+-----+---+-----+  
| id|  name|group_id| id|department|  
+---+-----+-----+---+-----+  
|  0|   Seif|      0|  0|  SICS/KTH|  
|  1|   Amir|      1|  1|      KTH|  
|  2|Sarunas|      1|  1|      KTH|  
+---+-----+-----+---+-----+
```

Joins Example - Outer

```
val joinExpression = person.col("group_id") === group.col("id")  
  
var joinType = "outer"  
  
person.join(group, joinExpression, joinType).show()
```

```
+---+-----+-----+---+-----+  
| id|  name|group_id| id|department|  
+---+-----+-----+---+-----+  
|  1|  Amir|      1|  1|      KTH|  
|  2|Sarunas|      1|  1|      KTH|  
|null|  null|   null|  2|      SICS|  
|  0|  Seif|      0|  0| SICS/KTH|  
+---+-----+-----+---+-----+
```

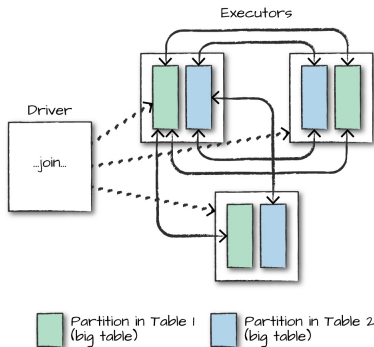


Joins Communication Strategies

- ▶ Two different communication ways during joins:
 - Shuffle join: big table to big table
 - Broadcast join: big table to small table

Shuffle Join

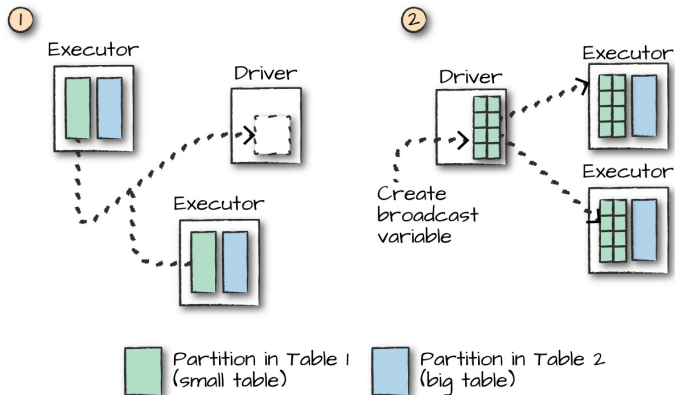
- ▶ Every node talks to **every other node**.
- ▶ They share data according to **which node** has a **certain key or set of keys**.



[M. Zaharia et al., Spark: The Definitive Guide, O'Reilly Media, 2018]

Broadcast Join

- When the table is **small enough** to **fit into the memory of a single worker node**.



[M. Zaharia et al., Spark: The Definitive Guide, O'Reilly Media, 2018]

SQL

- ▶ You can run **SQL queries** on views/tables via the method `sql` on the `SparkSession` object.

```
spark.sql("SELECT * from people_view").show()
```

```
+---+---+-----+
|age| id|  name|
+---+---+-----+
| 15| 12|Michael|
| 30| 15|  Andy|
| 19| 20| Justin|
| 12| 15|  Andy|
| 19| 20|   Jim|
| 12| 10|  Andy|
+---+---+-----+
```



Temporary View

- ▶ `createOrReplaceTempView` creates (or replaces) a lazily evaluated `view`.
- ▶ You can use it like a `table` in Spark SQL.

```
people.createOrReplaceTempView("people_view")  
  
val teenagersDF = spark.sql("SELECT name, age FROM people_view WHERE age BETWEEN 13 AND 19")
```

DataSet



Untyped API with DataFrame

- ▶ DataFrames elements are Rows, which are generic untyped JVM objects.
- ▶ Scala compiler cannot type check Spark SQL schemas in DataFrames.



Untyped API with DataFrame

- ▶ DataFrames elements are **Rows**, which are **generic untyped JVM objects**.
- ▶ Scala compiler **cannot type check** Spark SQL **schemas** in DataFrames.
- ▶ The following code **compiles**, but you get a **runtime exception**.
 - `id_num` is not in the DataFrame columns `[name, age, id]`

```
// people columns: ("name", "age", "id")  
val people = spark.read.format("json").load("people.json")  
  
people.filter("id_num < 20") // runtime exception
```



Why DataSet?

- ▶ Assume the following example

```
case class Person(name: String, age: BigInt, id: BigInt)
val peopleRDD = sc.parallelize(Array(Person("seif", 65, 0), Person("amir", 40, 1)))
val peopleDF = peopleRDD.toDF
```



Why DataSet?

- ▶ Assume the following example

```
case class Person(name: String, age: BigInt, id: BigInt)
val peopleRDD = sc.parallelize(Array(Person("seif", 65, 0), Person("amir", 40, 1)))
val peopleDF = peopleRDD.toDF
```

- ▶ Now, let's use `collect` to bring back it to the master.

```
val collectedPeople = peopleDF.collect()
// collectedPeople: Array[org.apache.spark.sql.Row]
```




Why DataSet?

- ▶ Assume the following example

```
case class Person(name: String, age: BigInt, id: BigInt)
val peopleRDD = sc.parallelize(Array(Person("seif", 65, 0), Person("amir", 40, 1)))
val peopleDF = peopleRDD.toDF
```

- ▶ Now, let's use `collect` to bring back it to the master.

```
val collectedPeople = peopleDF.collect()
// collectedPeople: Array[org.apache.spark.sql.Row]
```

- ▶ What is in `Row`?



Why DataSet?

- ▶ To be able to work with the collected values, we should **cast** the **Rows**.
 - How many **columns**?
 - What **types**?

```
// Person(name: Sting, age: BigInt, id: BigInt)

val collectedList = collectedPeople.map {
  row => (row(0).asInstanceOf[String], row(1).asInstanceOf[Int], row(2).asInstanceOf[Int])
}
```



Why DataSet?

- ▶ To be able to work with the collected values, we should **cast** the **Rows**.
 - How many **columns**?
 - What **types**?

```
// Person(name: Sting, age: BigInt, id: BigInt)

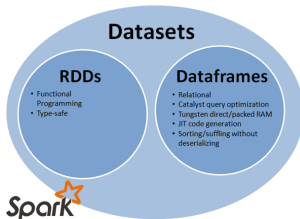
val collectedList = collectedPeople.map {
  row => (row(0).asInstanceOf[String], row(1).asInstanceOf[Int], row(2).asInstanceOf[Int])
}
```

- ▶ But, what if we cast the **types wrong**?
- ▶ Wouldn't it be nice if we could have both **Spark SQL optimizations** and **typesafety**?

DataSet

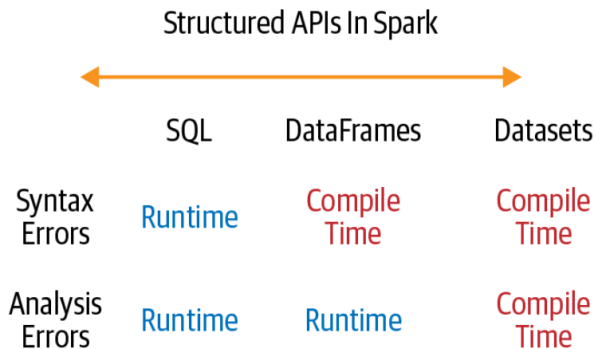
- ▶ **Datasets** can be thought of as **typed** distributed collections of data.
- ▶ **Dataset** API unifies the **DataFrame** and **RDD** APIs.
- ▶ You can consider a **DataFrame** as an alias for **Dataset [Row]**, where a **Row** is a **generic untyped JVM object**.

```
type DataFrame = Dataset[Row]
```



[<http://why-not-learn-something.blogspot.com/2016/07/apache-spark-rdd-vs-dataframe-vs-dataset.html>]

Structured APIs in Spark



[J.S. Damji et al., Learning Spark - Lightning-Fast Data Analytics]



Creating DataSets

- ▶ To convert a **sequence** or an **RDD** to a **Dataset**, we can use `toDS()`.
- ▶ You can call `as[SomeCaseClass]` to convert the **DataFrame** to a Dataset.

```
case class Person(name: String, age: BigInt, id: BigInt)
val personSeq = Seq(Person("Max", 33, 0), Person("Adam", 32, 1))
```



Creating DataSets

- ▶ To convert a **sequence** or an **RDD** to a **Dataset**, we can use `toDS()`.
- ▶ You can call `as[SomeCaseClass]` to convert the **DataFrame** to a Dataset.

```
case class Person(name: String, age: BigInt, id: BigInt)
val personSeq = Seq(Person("Max", 33, 0), Person("Adam", 32, 1))
```

```
val ds1 = sc.parallelize(personSeq).toDS
```



Creating DataSets

- ▶ To convert a **sequence** or an **RDD** to a **Dataset**, we can use `toDS()`.
- ▶ You can call `as[SomeCaseClass]` to convert the **DataFrame** to a Dataset.

```
case class Person(name: String, age: BigInt, id: BigInt)
val personSeq = Seq(Person("Max", 33, 0), Person("Adam", 32, 1))
```

```
val ds1 = sc.parallelize(personSeq).toDS
```

```
val ds2 = spark.read.format("json").load("people.json").as[Person]
```




DataSet Transformations

- ▶ Transformations on Datasets are the same as those that we had on DataFrames.
- ▶ Datasets allow us to specify more complex and strongly typed transformations.

```
case class Person(name: String, age: BigInt, id: BigInt)

val people = spark.read.format("json").load("people.json").as[Person]

people.filter(x => x.age < 40).show()

people.map(x => (x.name, x.age + 5, x.id)).show()
```





Questions

- ▶ What happens when **missing data** is dropped or filled in?



Possible Answers

- ▶ People and experiences can disappear from the dataset.



Possible Answers

- ▶ People and experiences can disappear from the dataset.
- ▶ Filling in with defaults (0, “unknown”) can create a false sense of completeness.



Possible Answers

- ▶ People and experiences can disappear from the dataset.
- ▶ Filling in with defaults (0, “unknown”) can create a false sense of completeness.
- ▶ Structural inequalities that caused the missingness get hidden.



Possible Answers

- ▶ People and experiences can disappear from the dataset.
- ▶ Filling in with defaults (0, “unknown”) can create a false sense of completeness.
- ▶ Structural inequalities that caused the missingness get hidden.
- ▶ Alternatives: keep missingness visible, trace why it is missing, or represent uncertainty instead of erasing it.



Alternatives

- ▶ Probabilistic databases (e.g., MayBMS, MystiQ)
 - Instead of filling missing values with a single guess, they represent uncertainty explicitly and let queries return probability distributions.

- ▶ **Probabilistic databases** (e.g., MayBMS, MystiQ)
 - Instead of filling missing values with a single guess, they represent uncertainty explicitly and let queries return probability distributions.

- ▶ **Data provenance systems** (e.g., DataHub, ProvSQL)
 - They attach metadata to each value showing its source and whether it's missing, inferred, or uncertain.



Questions

- ▶ What are the risks of forcing data into **strict structures** (schemas, categories, types)?



Possible Answers



Possible Answers

- ▶ Complex identities (e.g., gender, race, culture) get flattened into rigid boxes.



Possible Answers

- ▶ **Complex identities** (e.g., gender, race, culture) get flattened into **rigid boxes**.
- ▶ Values that don't fit the schema may be **deleted** or **misclassified**.



Possible Answers

- ▶ **Complex identities** (e.g., gender, race, culture) get flattened into **rigid boxes**.
- ▶ Values that don't fit the schema may be **deleted** or **misclassified**.
- ▶ Systems reflect the **worldview of the schema designer**, not the data subjects.



Possible Answers

- ▶ **Complex identities** (e.g., gender, race, culture) get flattened into **rigid boxes**.
- ▶ Values that don't fit the schema may be **deleted** or **misclassified**.
- ▶ Systems reflect the **worldview of the schema designer**, not the data subjects.
- ▶ **Alternatives:** allow flexible schemas, self-identification, multiple categories, or context-aware structures.



Alternatives

- ▶ **Document-oriented** NoSQL DBs (e.g., MongoDB, CouchDB)
 - They allow semi-structured data (different documents can have different fields), so not everything needs to fit a fixed schema.



Alternatives

- ▶ **Document-oriented** NoSQL DBs (e.g., MongoDB, CouchDB)
 - They allow semi-structured data (different documents can have different fields), so not everything needs to fit a fixed schema.
- ▶ **Graph databases** (e.g., Neo4j)
 - They don't require uniform rows/columns; nodes can have different properties, making it easier to represent irregular or relational data.

Summary



Summary

- ▶ RDD vs. DataFrame vs. DataSet
- ▶ Transformation and Actions
- ▶ Aggregation
- ▶ Join
- ▶ SQL queries



References

- ▶ M. Zaharia et al., “Spark: The Definitive Guide”, O'Reilly Media, 2018 - Chapters 4-11.
- ▶ M. Armbrust et al., “Spark SQL: Relational data processing in spark”, ACM SIGMOD, 2015.

Questions?