Spark DataFrames

Anurag Naga

Introduction

Creating DataFrames
Loading Data

Operations using DF

Projection
Ordering
Grouping
Joining

Spark DataFrames

Anurag Nagar

Big Data Class

Spark DataFrames

Anurag Naga

Introduction DataFrames

Creating DataFrame Loading Data

Operation using DF

Selection and Projection Ordering Grouping Joining

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data
- 2 Operations using DF
 - Selection and Projection
 - Ordering
 - Grouping
 - Joining

Introduction

Spark DataFrames

Anurag Naga

Introduction
DataFrames
Creating DataFrame
Loading Data

Operations
using DF
Selection and
Projection
Ordering

DataFrames are part of Spark SQL.

- Like RDDs, DataFrames (DF) are immutable, distributed, partitioned collection of data
- They have all the properties of RDDs, such as lazy evaluation, recovery through lineage graphs, etc.
- They contain specialized APIs for working with tabular data, and have named columns.

Name	Age	Height
String	Int	Double
String	Int	Double
String	Int	Double

String	Int	Double
String	Int	Double
String	Int	Double

DataFrame

Spark DataFrames

Anurag Nagai

Introduction

DataFrames

Creating DataFram

Creating DataFrames Loading Data

Operations using DF

Selection and Projection Ordering Grouping Joining

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data
- 2 Operations using DF
 - Selection and Projection
 - Ordering
 - Grouping
 - Joining

Creating DataFrames

Spark DataFrames

Anurag Naga

ntroduction

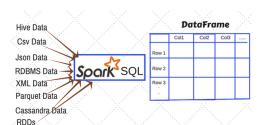
DataFrames

Creating DataFrames

Loading Data

Operations
using DF
Selection and
Projection
Ordering
Grouping
Joining

- DataFrames are well suited for large structured or semi-structured data.
- Data can be loaded easily from a wide variety of sources
- DF contain named columns, and a list of tuples



Spark DataFrames

Allulag Ivaga

Introduction DataFrames Creating DataFrame Loading Data

Operations using DF

Selection and Projection Ordering Grouping Joining

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data
- 2 Operations using DF
 - Selection and Projection
 - Ordering
 - Grouping
 - Joining

Loading Data into DataFrames

Spark DataFrames

Anurag Naga

Introduction
DataFrames
Creating DataFrames
Loading Data

Operations using DF

Selection and Projection Ordering Grouping Joining **spark.read** is the starting point to read data into DF. More details can be found at this link.

■ To read a simple CSV file with header

Loading Data into DataFrames

Spark DataFrames

Anurag Naga

Introduction
DataFrames
Creating DataFrames
Loading Data

Operations using DF Selection and Projection

Selection and Projection Ordering Grouping Joining **spark.read** is the starting point to read data into DF. More details can be found at this link.

■ To read a simple CSV file with header

■ To see schema

```
df.printSchema()\\
```

Loading Data into DataFrames

Spark DataFrames

Anurag Naga

Introduction
DataFrames
Creating DataFrames
Loading Data

Operations using DF Selection and Projection Ordering **spark.read** is the starting point to read data into DF. More details can be found at this link.

■ To read a simple CSV file with header

■ To see schema

```
df.printSchema()
```

To see first 10 rows

```
df.take(10)
```

Spark DataFrames

Anurag Naga

Introduction

DataFrames

Creating DataFrames

Loading Data

Operations using DF Selection and Projection

Selection and Projection Ordering Grouping

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data

2 Operations using DF

- Selection and Projection
- Ordering
- Grouping
- Joining

Selection and Projection

Spark DataFrames

Anurag Naga

Introduction

DataFrames

Creating DataFrames

Loading Data

Operations using DF Selection and

Projection
Ordering
Grouping
Joining

To extract few columns

```
filtered = df. select (["column1", "column2"])
```

Selection and Projection

Spark DataFrames

Anurag Naga

Introduction
DataFrames
Creating DataFrames
Loading Data

Operations using DF Selection and Projection Ordering Grouping ■ To extract few columns

```
\mathsf{filtered} = \mathsf{df.}\,\mathsf{select}\,\big([\,\mathsf{"column1"},\,\mathsf{"column2"}]\big)
```

■ To filter data with conditions:

Spark **DataFrames**

Ordering

- DataFrames
- Creating DataFrames
- Loading Data

2 Operations using DF

- Selection and Projection
- Ordering
- Grouping
- Joining

Ordering

Spark DataFrames

Anurag Naga

Introduction

DataFrames

Creating DataFrames

Loading Data

Operations using DF Selection and Projection

Ordering
Grouping

■ To order by a column

Spark DataFrames

Anurag Naga

Introduction DataFrames Creating DataFrames

Operation using DF

Selection and Projection

Grouping

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data

2 Operations using DF

- Selection and Projection
- Ordering
- Grouping
- Joining

Grouping Data

Spark DataFrames

Anurag Naga

Introduction

DataFrames

Creating DataFrame

Operation

Selection and

Grouping

■ To group by a column and get count of groups:

```
\mathsf{df}.\mathsf{groupBy}(\texttt{"age"}).\mathsf{count}()
```

Grouping Data

Spark DataFrames

Allurag Ivaga

Introduction DataFrames Creating DataFrames Loading Data

Operations using DF

Selection and Projection

Grouping

■ To group by a column and get count of groups:

```
df.groupBy("age").count()
```

 To group by a column and show average of another column by group

```
df.groupBy("department").avg("salary")
```

Grouping Data

Spark **DataFrames**

Grouping

■ To group by a column and get count of groups:

```
df.groupBy("age").count()
```

To group by a column and show average of another column by group

```
df.groupBy("department").avg("salary")
```

To find other stats

```
df.groupBy("department")
    .agg(sum("salary").alias("sum_salary"),
         avg("salary"). alias ("avg_salary"),
         sum("bonus").alias("sum bonus"),
         max("bonus").alias("max_bonus"))
```

Spark DataFrames

Anurag Naga

Introduction DataFrames Creating DataFrame

Operations using DF

Selection and Projection Ordering Grouping Joining

1 Introduction

- DataFrames
- Creating DataFrames
- Loading Data

2 Operations using DF

- Selection and Projection
- Ordering
- Grouping
- Joining

Joining Data

Spark DataFrames

Anurag Naga

Introduction DataFrames Creating DataFrames

Operations

Selection an Projection Ordering Grouping To join two DF

```
df = left.join(right, left.name == right.name, "inner")
```

Joining Data

Spark DataFrames

Anurag Naga

Introduction

DataFrames

Creating DataFrames

Loading Data

Operations using DF

Selection and Projection Ordering Grouping Joining ■ To join two DF

```
\mathsf{df} \, = \, \mathsf{left} \, . \, \mathsf{join} \, \big( \, \mathsf{right} \, , \, \, \, \, \mathsf{left} \, . \, \mathsf{name} = = \, \mathsf{right.name}, \, \, "\mathsf{inner"} \big)
```

■ To do left/right outer join

the last parameter can be inner, outer, leftOuter, rightOuter