

CSCI 443 LECTURE 5: MUCH MORE ABOUT SPARK

Professor David Harrison

TODAY

- Spark Transformations and Actions
- Spark RDDs
- Spark ...

Spark



HOMework 1

Due Tonight.

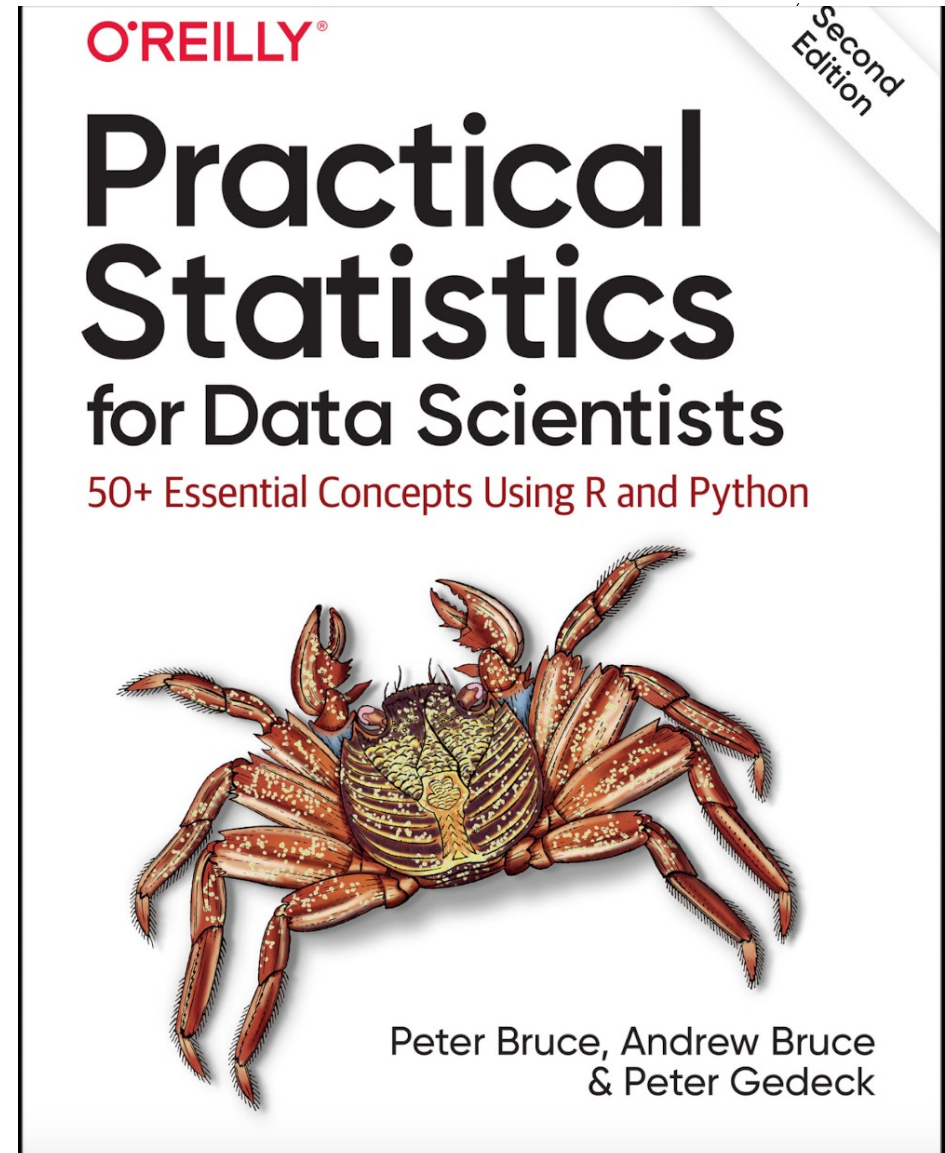
11 pm.

Submission:

- Submit archived Databricks Notebook to Blackboard.
- NOTE: Submission only needs to be the notebook. No README is necessary.

READING!

- Read Chapter 1: Exploratory Data Analysis.





OFFICE HOURS

Due to scheduling conflict, office hours updated

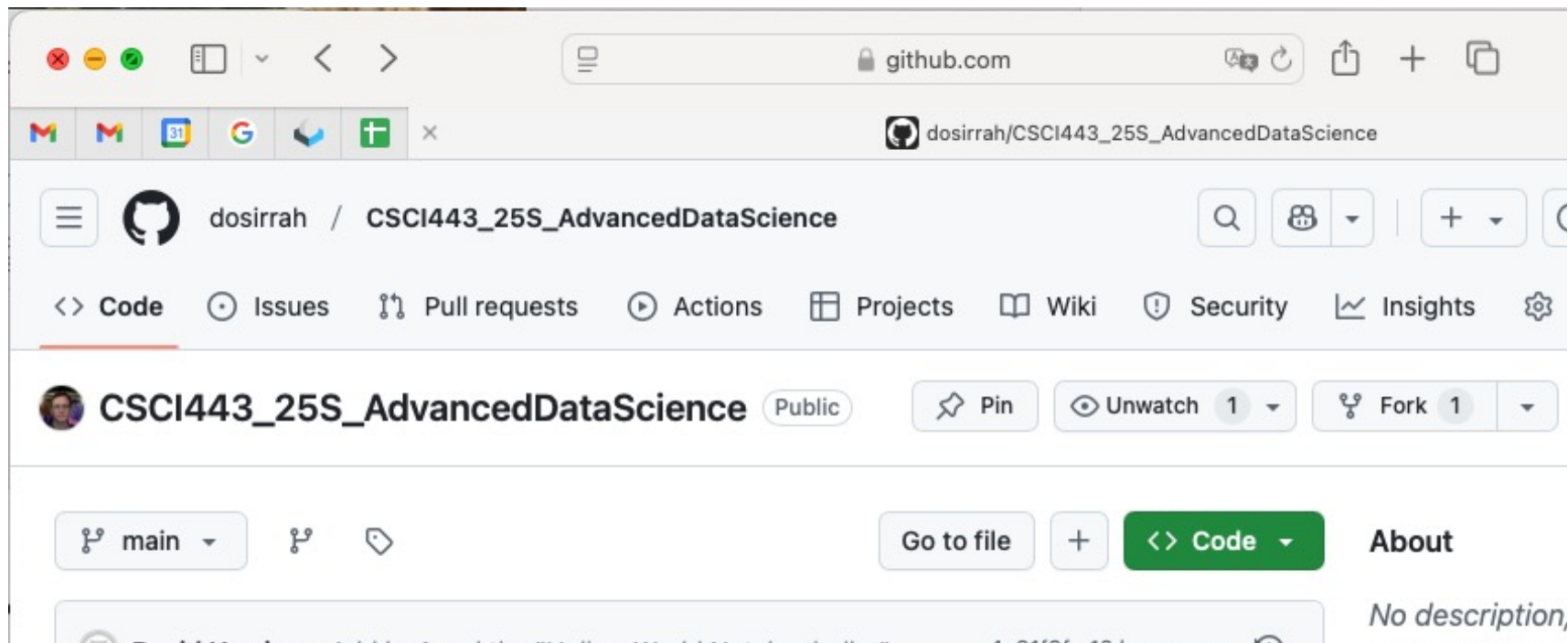
Monday	1:00-2:00 PM
Tuesday	4-5 PM

GITHUB

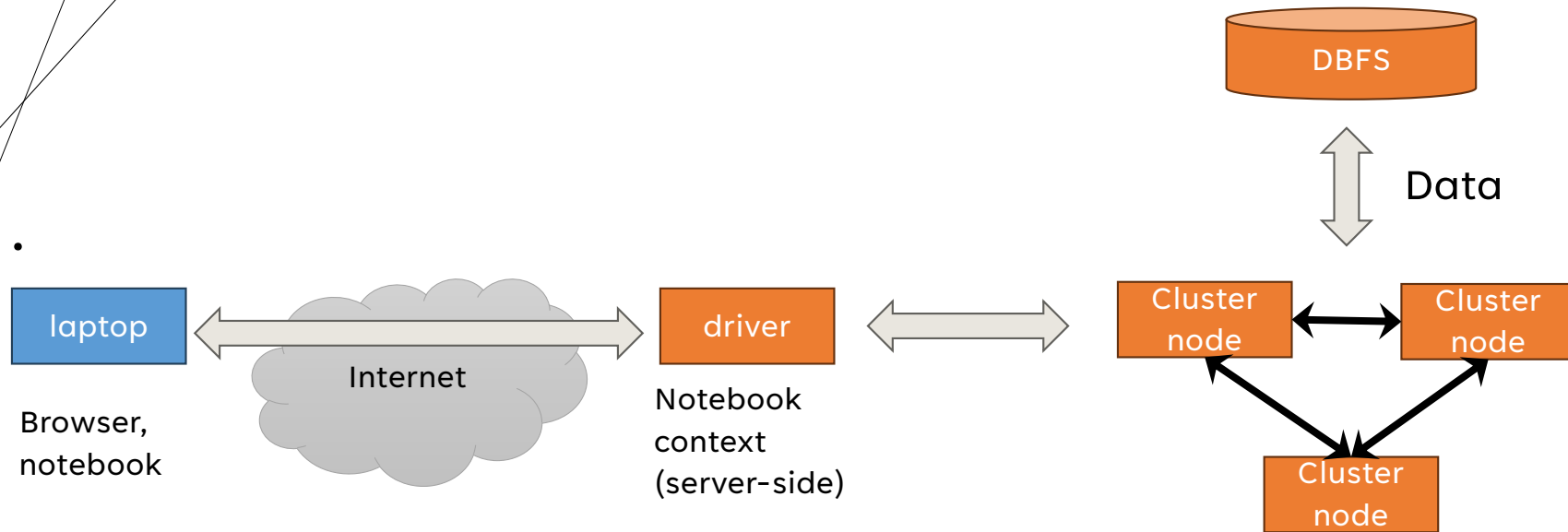
Lecture slides and examples have been committed to GitHub for lectures 1 through 4.

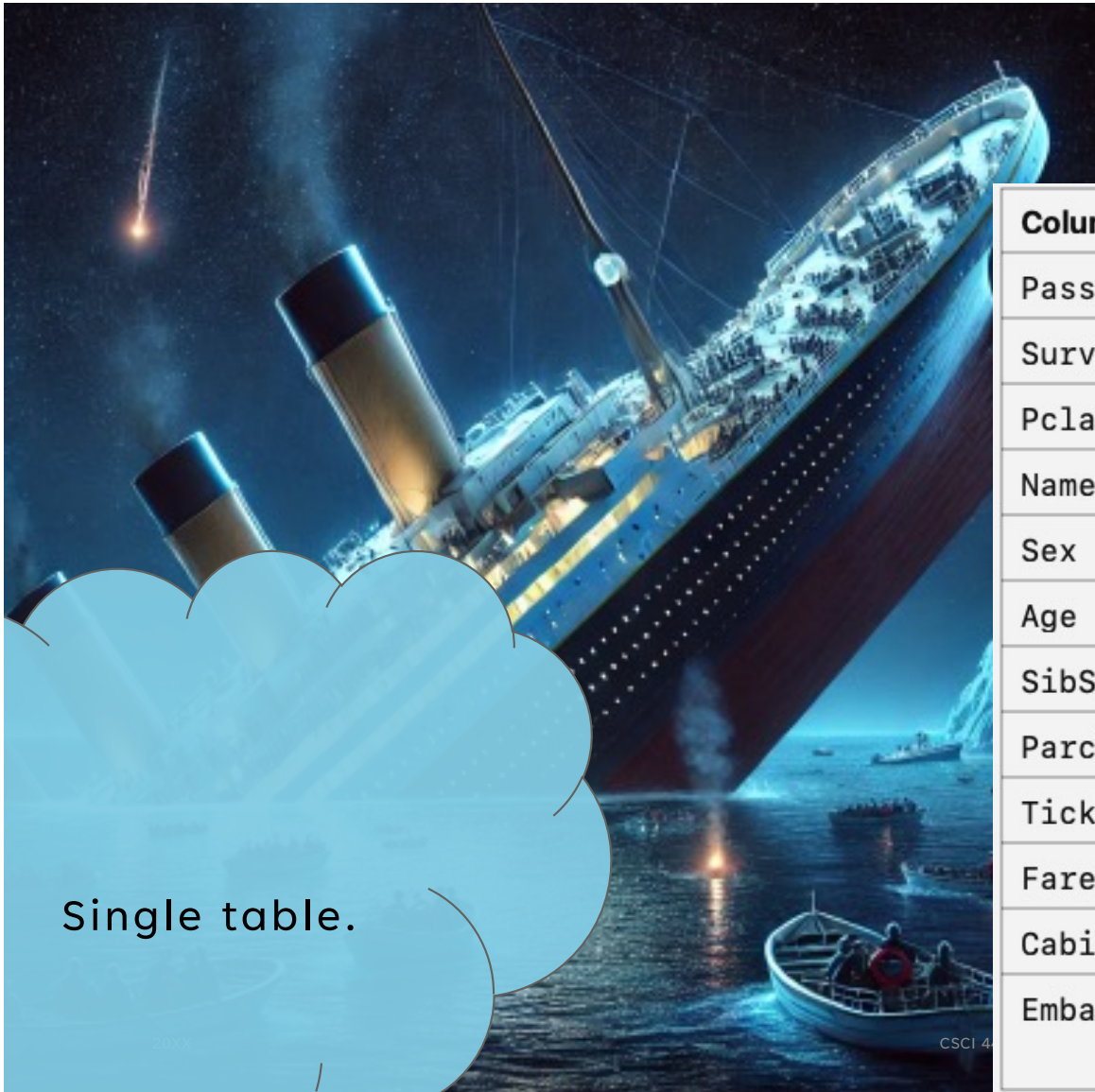
The project is at

https://github.com/dosirrah/CSCI443_25S_AdvancedDataScience



SPARK ARCHITECTURE (1000 FT)





TITANIC DATASET

Column Name	Description
PassengerId	Unique identifier for each passenger
Survived	Survival (0 = No, 1 = Yes)
Pclass	Passenger class (1 = 1st, 2 = 2nd, 3 = 3rd)
Name	Passenger's full name
Sex	Passenger's sex (male/female)
Age	Passenger's age (in years)
SibSp	Number of siblings/spouses aboard the Titanic
Parch	Number of parents/children aboard the Titanic
Ticket	Ticket number
Fare	Fare price paid for the ticket
Cabin	Cabin number (if assigned)
Embarked	Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)



Small table.

TITANIC DATASET

▶ ▼ ✓ Just now (1s)

```
df.count()
```

▶ (2) Spark Jobs

Out [14]: 891



Small table.

891 rows.

TITANIC DATASET

▶ ▼ ✓ Just now (1s)

```
df.count()
```

▶ (2) Spark Jobs

Out [14]: 891

TITANIC DATASET

Today we avoid using
Pandas or
pyspark.pandas.

Only Spark
DataFrames

By default Spark
does not infer a
schema.

▶ 07:35 PM (2s)

```
import pyspark.pandas as ps

file_path = "dbfs:/FileStore/shared_uploads/daharri6@olemiss.edu"

# Load a Spark DataFrame
df = spark.read.option("header", "true").csv(file_path)

# Print out the first 10 rows
df.show(10)
```

▶ (2) Spark Jobs

```
▼ df: pyspark.sql.dataframe.DataFrame
  PassengerId: string
    Survived: string
      Pclass: string
        Name: string
          Sex: string
            Age: string
              SibSp: string
                Parch: string
                  Ticket: string
                    Fare: string
                      Cabin: string
                        Embarked: string
```


TITANIC DATASET

Today we avoid using
Pandas or
pyspark.pandas.

Only Spark
DataFrames

By default Spark
does not infer a
schema.

```
df = spark.read.option("header", "true").option("inferSchema", "true").csv(file_path)
df.printSchema()
```

▸ (2) Spark Jobs

▸  df: pyspark.sql.dataframe.DataFrame = [PassengerId: integer, Survived: integer ... 10 more fields]

root

```
|-- PassengerId: integer (nullable = true)
|-- Survived: integer (nullable = true)
|-- Pclass: integer (nullable = true)
|-- Name: string (nullable = true)
|-- Sex: string (nullable = true)
|-- Age: double (nullable = true)
|-- SibSp: integer (nullable = true)
|-- Parch: integer (nullable = true)
|-- Ticket: string (nullable = true)
|-- Fare: double (nullable = true)
|-- Cabin: string (nullable = true)
|-- Embarked: string (nullable = true)
```

TITANIC DATASET

```
df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns]).show()
```

Every Spark operation is either a transformation or an operation.

Type	Description	Execution time
Transformation	Defines an operation that is added to the plan. No need to specify on what data the plan will be executed while building the plan.	Deferred
Action	Executes all transformations in the the plan on the specified data.	Immediate

TITANIC DATASET

```
df.select(col("Cabin").isNull())
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column <i>c</i> when executed
isNull	deferred	Transformation to return true for each field that is null.
df.select	?	Selects a column on Dataframe <i>df</i>

What should
this return?



?

TITANIC DATASET

```
df.select(col("Cabin").isNull())
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column <i>c</i> when executed
isNull	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe <i>df</i>



```
DataFrame[(Cabin IS NULL): boolean]
```

A new
DataFrame?

TITANIC DATASET

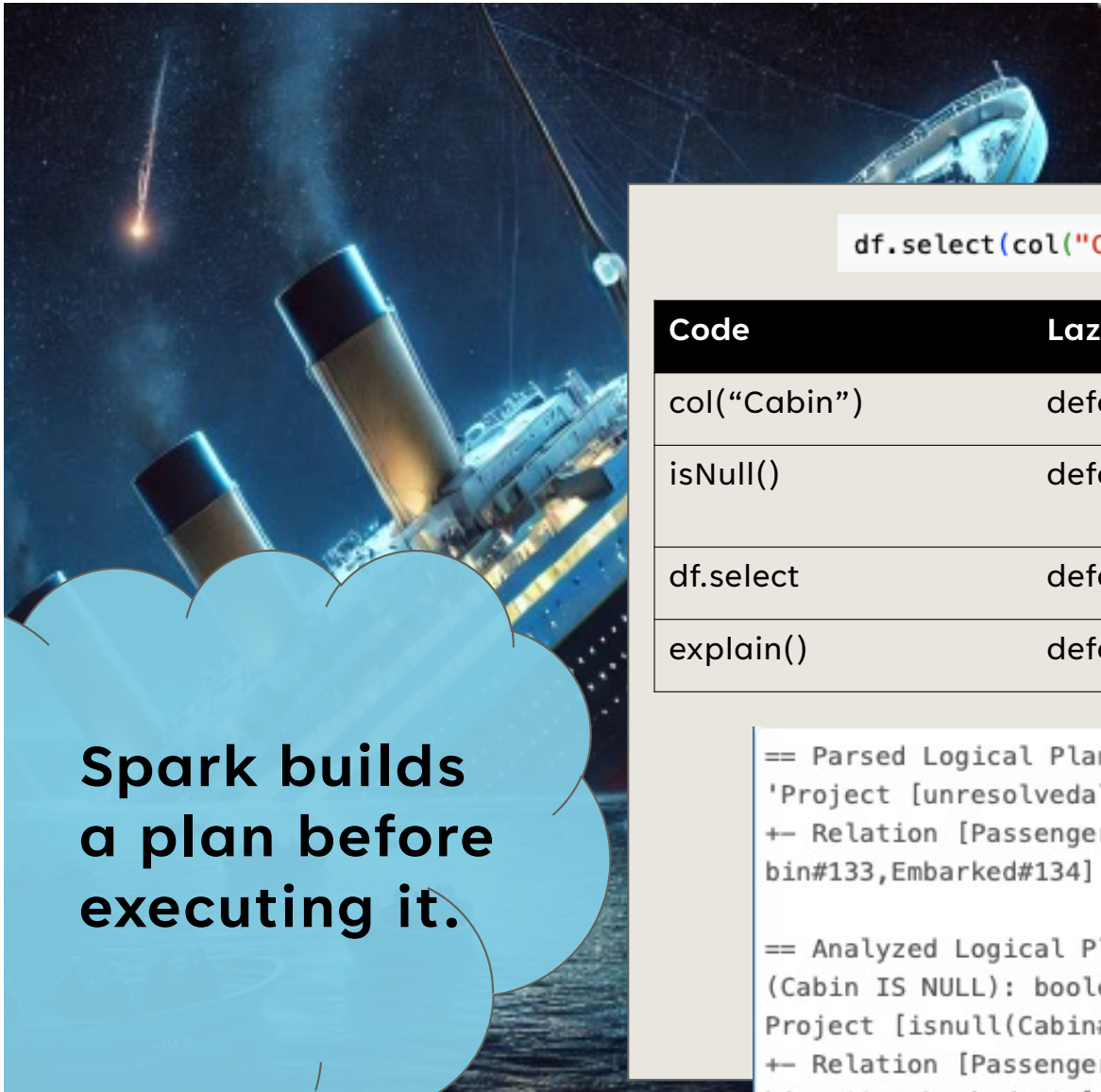
```
df.select(col("Cabin").isNull()).explain(extended=True)
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe <i>df</i>
explain()	deferred	<i>Describes the plan.</i>

A new
DataFrame?
Use explain()

```
== Parsed Logical Plan ==
'Project [unresolvedalias(isnull('Cabin), Some(org.apache.spark
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
bin#133,Embarked#134] csv

== Analyzed Logical Plan ==
(Cabin IS NULL): boolean
Project [isnull(Cabin#133) AS (Cabin IS NULL)#1686]
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
```



TITANIC DATASET

```
df.select(col("Cabin").isNull()).explain(extended=True)
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe df
explain()	deferred	<i>Describes the plan.</i>



```
== Parsed Logical Plan ==
'Project [unresolvedalias(isnull('Cabin), Some(org.apache.spark
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
bin#133,Embarked#134] csv

== Analyzed Logical Plan ==
(Cabin IS NULL): boolean
Project [isnull(Cabin#133) AS (Cabin IS NULL)#1686]
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
```

Spark builds
a plan before
executing it.

TITANIC DATASET

```
df.select(col("Cabin").isNull()).explain(extended=True)
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe <i>df</i>
explain()	deferred	<i>Describes the plan.</i>

We have yet to execute anything on the cluster.

```
== Parsed Logical Plan ==
'Project [unresolvedalias(isnull('Cabin), Some(org.apache.spark
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
bin#133,Embarked#134] csv

== Analyzed Logical Plan ==
(Cabin IS NULL): boolean
Project [isnull(Cabin#133) AS (Cabin IS NULL)#1686]
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
```


TITANIC DATASET

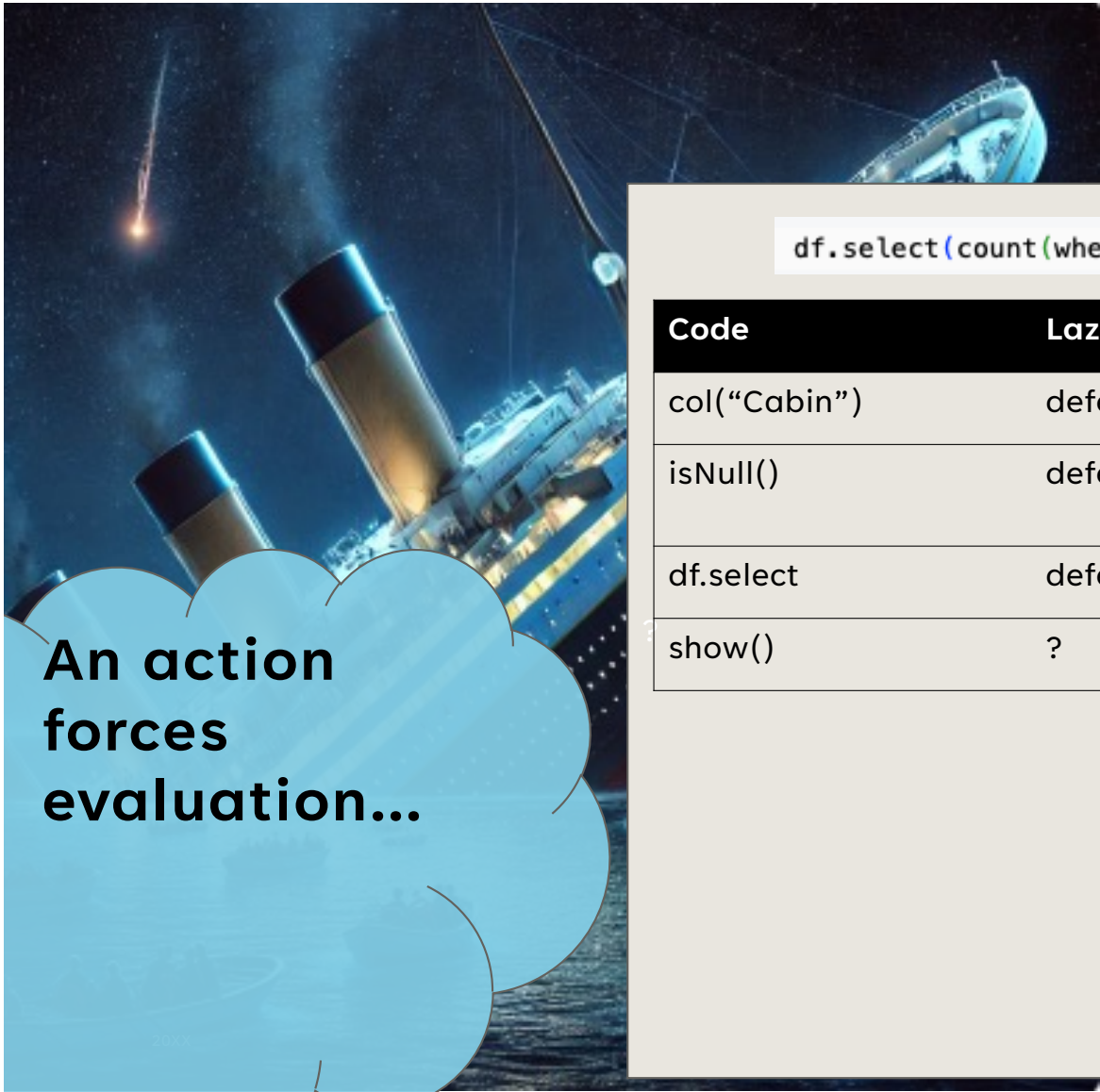
```
df.select(col("Cabin").isNull()).explain(extended=True)
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe <i>df</i>
explain()	deferred	<i>Describes the plan.</i>

Execution occurs when we introduce an *action*.

```
== Parsed Logical Plan ==
'Project [unresolvedalias(isnull('Cabin), Some(org.apache.spark
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
bin#133,Embarked#134] csv

== Analyzed Logical Plan ==
(Cabin IS NULL): boolean
Project [isnull(Cabin#133) AS (Cabin IS NULL)#1686]
+- Relation [PassengerId#123,Survived#124,Pclass#125,Name#126,S
```



An action
forces
evaluation...

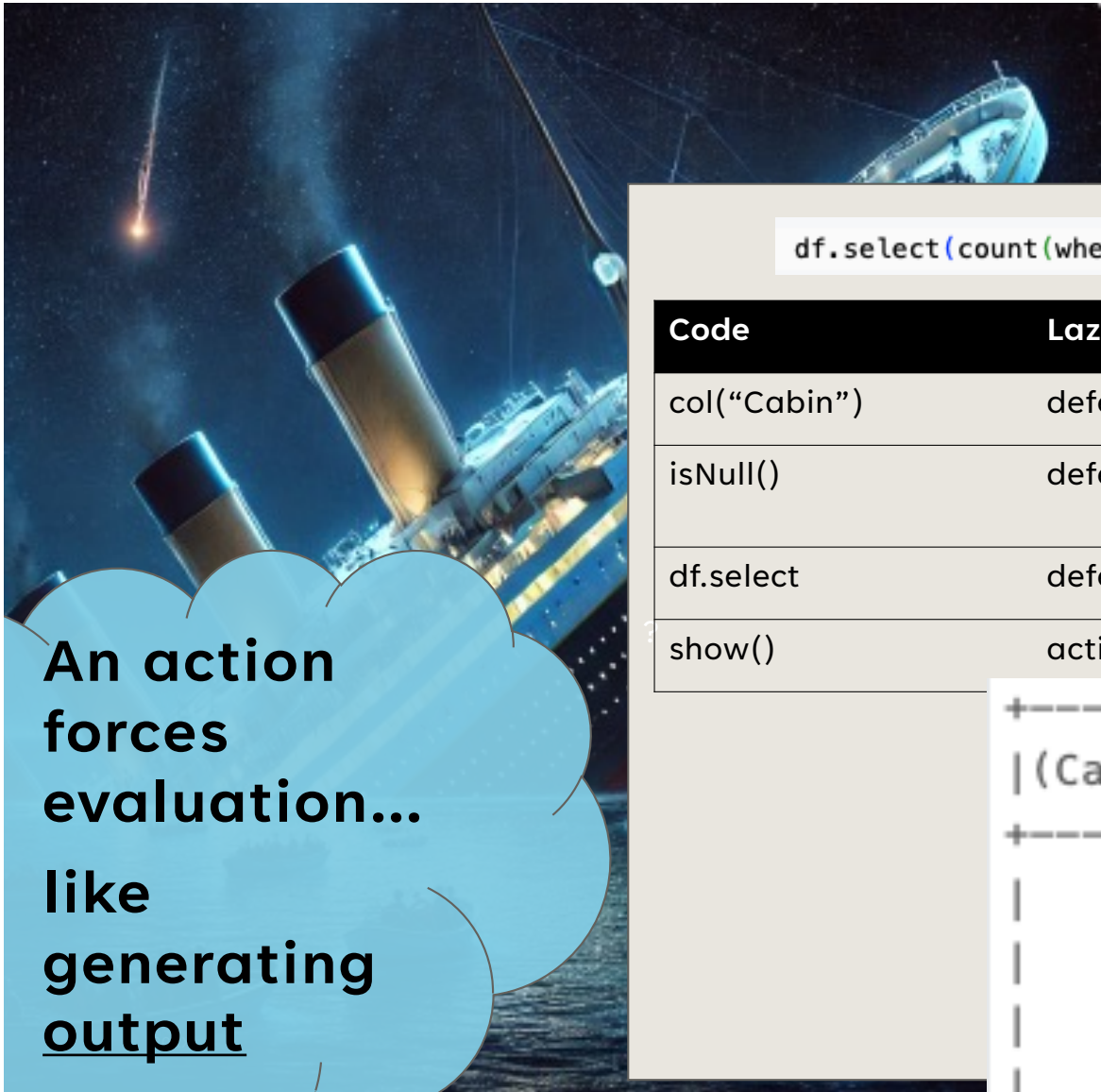
TITANIC DATASET

```
df.select(count(when(col("Cabin").isNull(), "Cabin"))).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe df
show()	?	Show us the result.



?



An action
forces
evaluation...
like
generating
output

TITANIC DATASET

```
df.select(count(when(col("Cabin").isNull(), "Cabin"))).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull()	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates plan with Dataframe df
show()	action	Show us the result.

↓

```
+-----+
| (Cabin IS NULL) |
+-----+
|                |
|                |
|                |
|                |
|                |
```

TITANIC DATASET

Generating output
forced execution.

The output is a
DataFrame
containing the
evaluation of `isNull`
for each cell in the
“Cabin” column.

```
df.select(col("Cabin").isNull()).show()
```

► (1) Spark Jobs

```
+-----+  
| (Cabin IS NULL) |  
+-----+  
|         true    |  
|        false    |  
|         true    |  
|        false    |  
|         true    |  
|         true    |  
|        false    |  
|         true    |  
|         true    |  
|         true    |  
|        false    |  
|        false    |  
|         true    |
```

TITANIC DATASET

```
df.select(col("Cabin").isNull())
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column <i>c</i> when executed
isNull	deferred	Transformation to return true for each field that is null.
df.select	deferred	Associates the plan with Dataframe <i>df</i>

What if I want to know how many nulls occurred in the "Cabin" column?

TITANIC DATASET

```
df.select(count(col("Cabin").isNull())).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column <i>c</i> when executed
isNull	deferred	Transformation to return true for each field that is null.
count	deferred	Add counting transformation to plan.
df.select	deferred	Associates the plan with Dataframe <i>df</i>
show	action	Show the result

```
+-----+
|count((Cabin IS NULL))|
+-----+
|                        891|
+-----+
```

What if I want to know how many nulls occurred in the "Cabin" column?

count()

TITANIC DATASET

```
df.select(count(col("Cabin").isNull())).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull	deferred	Transformation to return true for each field that is null.
count	deferred	Add counting transformation to plan.
df.select	deferred	Associates the plan with Dataframe <i>df</i>
show	action	Show the result

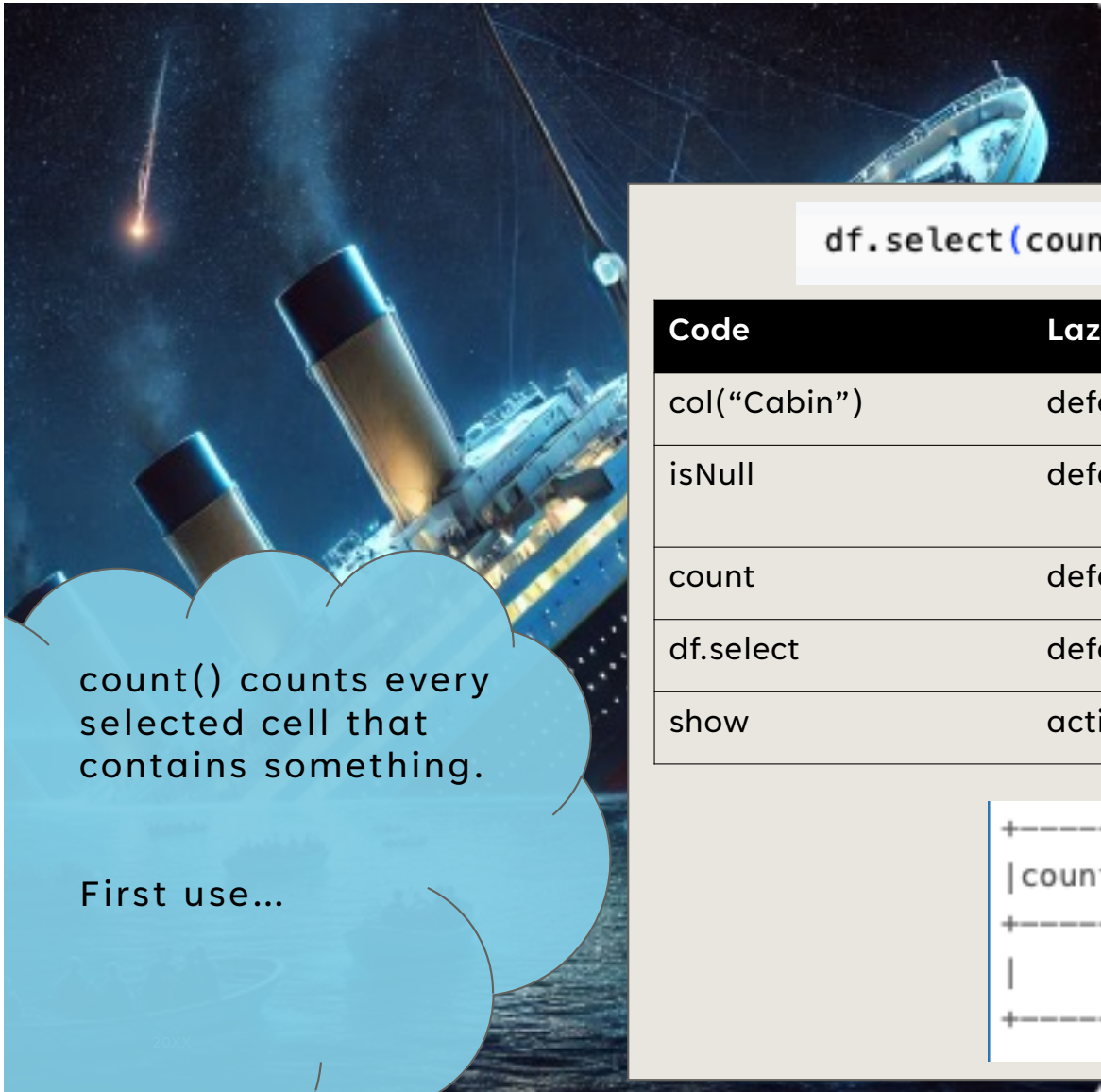
What if I want to know how many nulls occurred in the "Cabin" column?

count()

```
+-----+  
|count((Cabin IS NULL))|  
+-----+  
|                        |  
+-----+  
|                        |  
+-----+
```

891

X No



TITANIC DATASET

```
df.select(count(col("Cabin").isNull())).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull	deferred	Transformation to return true for each field that is null.
count	deferred	Add counting transformation to plan.
df.select	deferred	Associates the plan with Dataframe df
show	action	Show the result



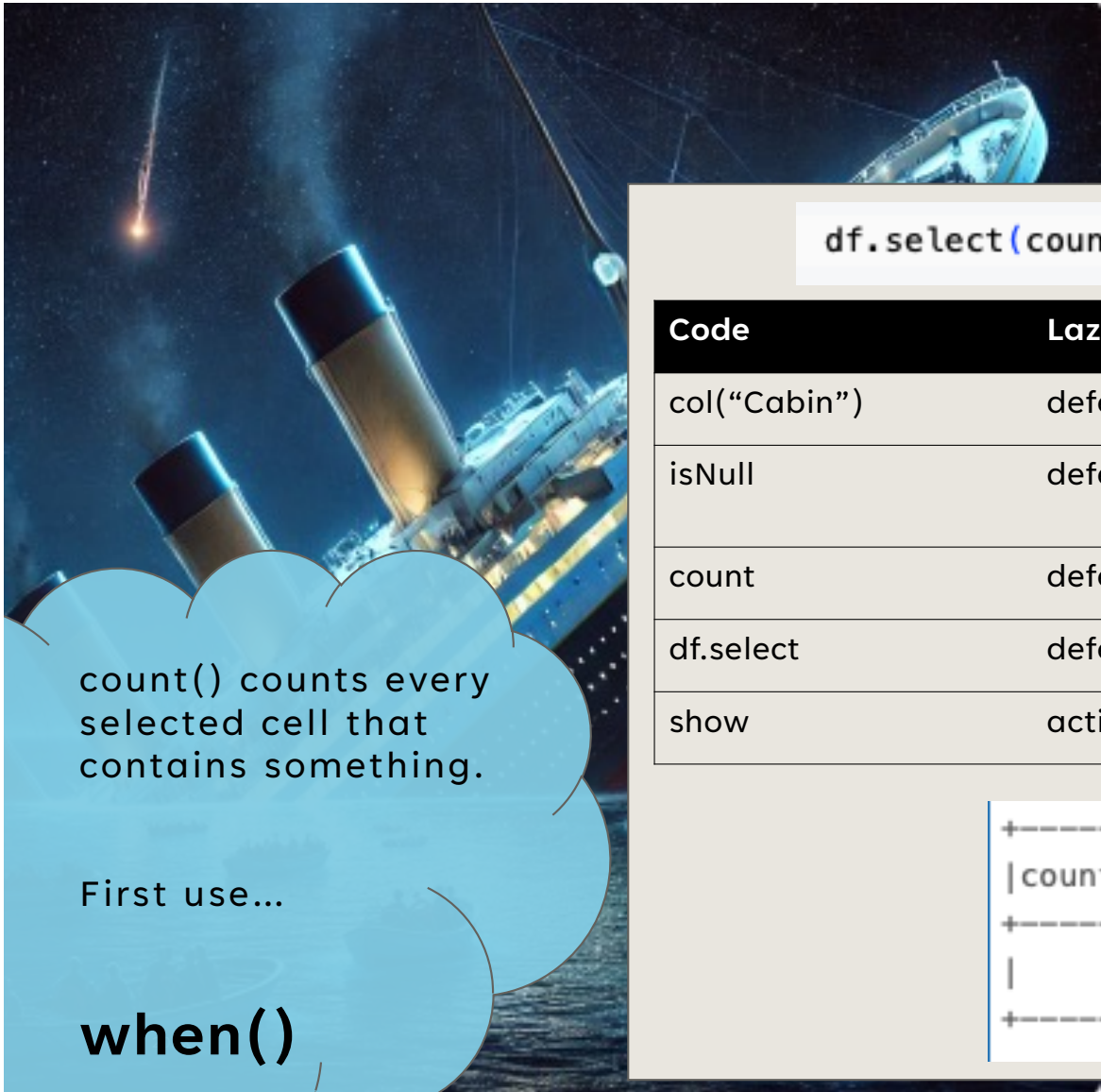
```
+-----+
|count((Cabin IS NULL))|
+-----+
|                        |
+-----+
|                        |
+-----+
```

891

X No

count() counts every selected cell that contains something.

First use...



TITANIC DATASET

```
df.select(count(col("Cabin").isNull())).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull	deferred	Transformation to return true for each field that is null.
count	deferred	Add counting transformation to plan.
df.select	deferred	Associates the plan with Dataframe df
show	action	Show the result



```
+-----+
|count((Cabin IS NULL))|
+-----+
|                        |
+-----+
```

891

X No

count() counts every selected cell that contains something.

First use...

when()

TITANIC DATASET

```
df.select(count(when(col("Cabin").isNull(), 1))).show()
```

Code	Lazy?	Description
col("Cabin")	deferred	Will select column c when executed
isNull	deferred	Transformation to return true for each field that is null.
when	deferred	Replace nulls with ones.
count	deferred	Add counting transformation to plan.
df.select	deferred	Associates the plan with Dataframe df
show	action	Show the result

```
+-----+
|count(CASE WHEN (Cabin IS NULL) THEN Cabin END)|
+-----+
|                                                    687|
+-----+
```

count() counts every selected cell that contains something.

First use...

when()

TITANIC DATASET

Last time we did some Exploratory Data Analysis, we looked for missing ages.

Let search for any other missing data.

```
from pyspark.sql.functions import col, isnull, when, count
```

```
df.select([count.when(col(c).isnull(), c)).alias(c) for c in df.columns]).show()
```

► (2) Spark Jobs

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	0	0	0	0	177	0	0	0	0	687	2

TRANSFORMATIONS

Added to the plan.

Transformation	Description	Example
<code>select()</code>	Select specific columns	<code>df.select("Name", "Age")</code>
<code>filter()/where()</code>	Filter rows based on a condition	<code>df.filter(df.Age > 30)</code>
<code>withColumn()</code>	Add or modify a column	<code>df.withColumn("AgePlusOne", df.Age + 1)</code>
<code>drop()</code>	Remove a column	<code>df.drop("Ticket")</code>
<code>orderBy()</code>	Sort rows by a column	<code>df.orderBy("Fare", ascending=False)</code>
<code>limit()</code>	Take the first N rows (still lazy)	<code>df.limit(10)</code>
<code>distinct()</code>	Remove duplicate rows	<code>df.select("Pclass").distinct()</code>

ACTIONS

Execute the plan to generate output.

Action	Description	Example
<code>show()</code>	Displays results	<code>df.show(5)</code>
<code>collect()</code>	Brings all data to the driver (⚠ not recommended for large data)	<code>df.collect()</code>
<code>count()</code>	Returns the total number of rows	<code>df.count()</code>
<code>first()</code>	Returns the first row	<code>df.first()</code>
<code>take(n)</code>	Returns the first n rows	<code>df.take(5)</code>
<code>describe()</code>	Computes summary statistics	<code>df.describe().show()</code>
<code>summary()</code>	More detailed statistics than <code>describe()</code>	<code>df.summary().show()</code>

RESILIENT DISTRIBUTED DATASETS (RDD)

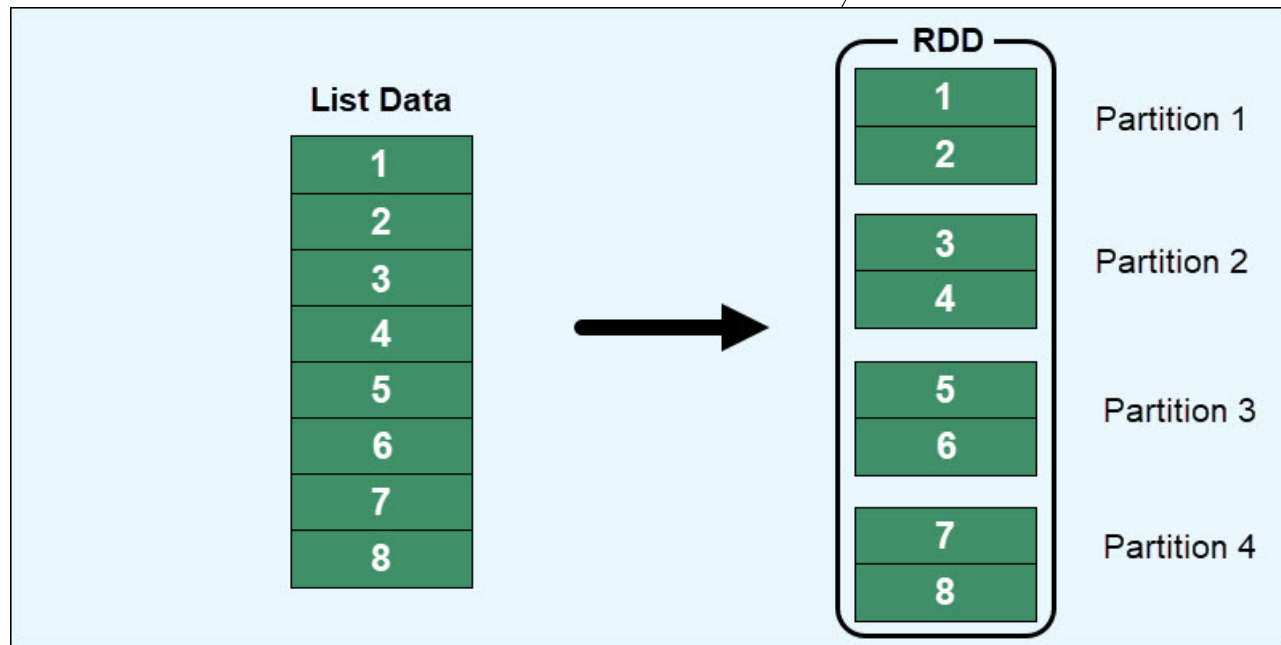
Each DataFrame has an underlying RDD.

Each Row in the RDD holds data for a record.

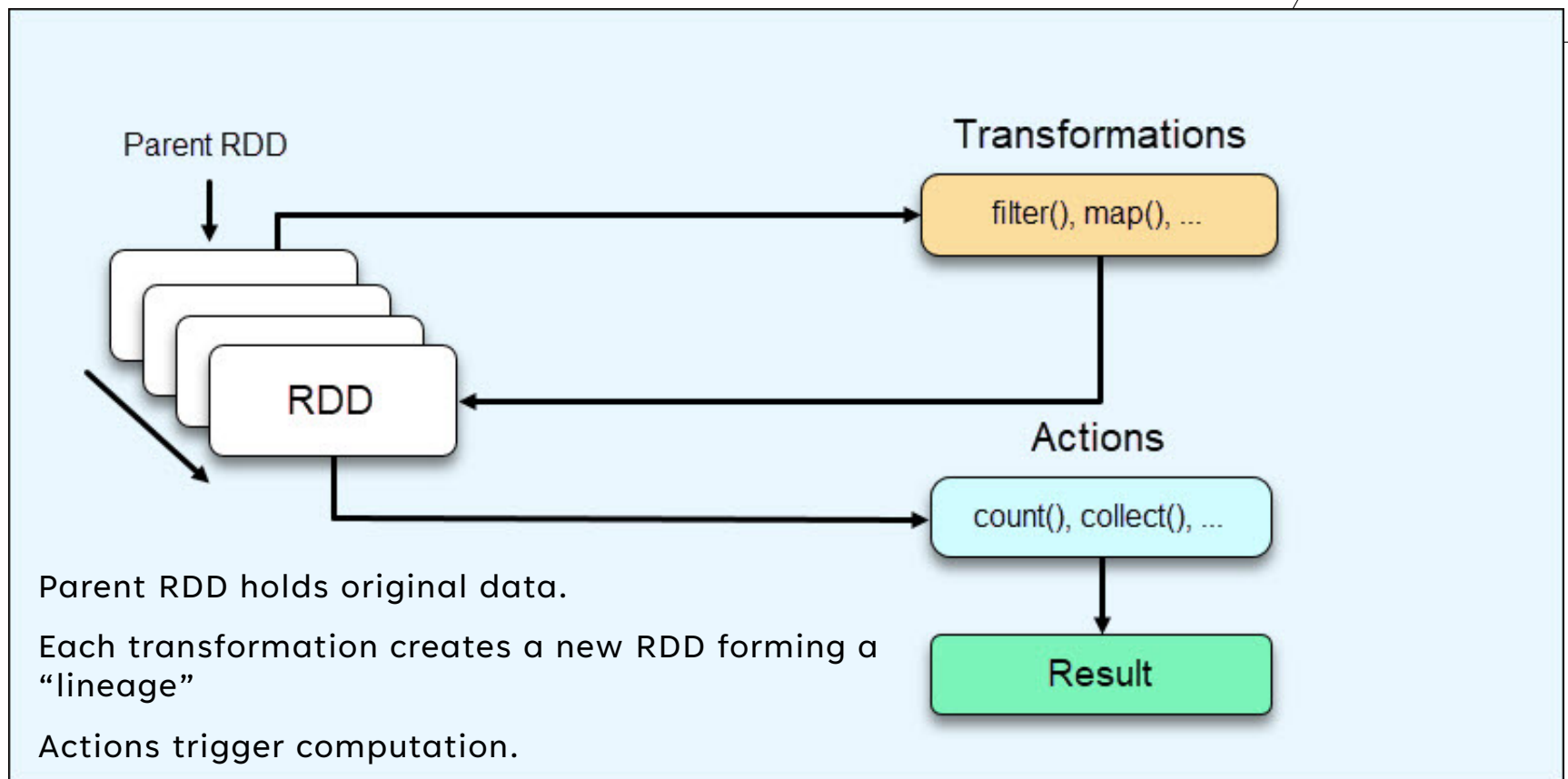
A Row is like an array where each field corresponds to a column in the DataFrame.

RDD partitioned across the cluster.

- Set of rows in each partition.



SEQUENCE OF TRANSFORMATIONS





THANK YOU

David Harrison

Harrison@cs.olemiss.edu