

INTRODUCTION TO

---

**APACHE SPARK**

---

# WHAT IS APACHE SPARK?

- ▶ <https://databricks.com/spark/about>
- ▶ <http://spark.apache.org>

---

# RDD

- ▶ Immutable Distributed Collection of Objects.
- ▶ Ways to Create a RDD :
  - ▶ Use external Dataset
    - ▶ `lines = sc.textFile("filename")`
  - ▶ Parallelize a collection
    - ▶ `lines = sc.parallelize([1,2,3,4,5,6])`

---

# TRANSFORMATION

- ▶ Always return new RDD's.
- ▶ Lazy Evaluation
- ▶ <http://spark.apache.org/docs/latest/programming-guide.html#transformations>
- ▶ `x = sc.parallelize([1,2,3,4,5])`
- ▶ `x10 = x.map(lambda n : n*10)`

*Table 3-2. Basic RDD transformations on an RDD containing {1, 2, 3, 3}*

Function name	Purpose	Example	Result
<code>map()</code>	Apply a function to each element in the RDD and return an RDD of the result.	<code>rdd.map(x =&gt; x + 1)</code>	<code>{2, 3, 4, 4}</code>
<code>flatMap()</code>	Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words.	<code>rdd.flatMap(x =&gt; x.to(3))</code>	<code>{1, 2, 3, 2, 3, 3, 3}</code>
<code>filter()</code>	Return an RDD consisting of only elements that pass the condition passed to <code>filter()</code> .	<code>rdd.filter(x =&gt; x != 1)</code>	<code>{2, 3, 3}</code>
<code>distinct()</code>	Remove duplicates.	<code>rdd.distinct()</code>	<code>{1, 2, 3}</code>
<code>sample(withReplacement, fraction, [seed])</code>	Sample an RDD, with or without replacement.	<code>rdd.sample(false, 0.5)</code>	Nondeterministic

---

# ACTIONS

- ▶ Actual Computation happens.
- ▶ <http://spark.apache.org/docs/latest/programming-guide.html#actions>
- ▶ `sum = x10.reduce(lambda n,m : n+m)`

*Table 3-4. Basic actions on an RDD containing {1, 2, 3, 3}*

Function name	Purpose	Example	Result
<code>collect()</code>	Return all elements from the RDD.	<code>rdd.collect()</code>	{1, 2, 3, 3}
<code>count()</code>	Number of elements in the RDD.	<code>rdd.count()</code>	4
<code>countByValue()</code>	Number of times each element occurs in the RDD.	<code>rdd.countByValue()</code>	{(1, 1), (2, 1), (3, 2)}

Function name	Purpose	Example	Result
<code>take(num)</code>	Return num elements from the RDD.	<code>rdd.take(2)</code>	<code>{1, 2}</code>
<code>top(num)</code>	Return the top num elements the RDD.	<code>rdd.top(2)</code>	<code>{3, 3}</code>
<code>takeOrdered(num)(ordering)</code>	Return num elements based on provided ordering.	<code>rdd.takeOrdered(2)(myOrdering)</code>	<code>{3, 3}</code>
<code>takeSample(withReplacement, num, [seed])</code>	Return num elements at random.	<code>rdd.takeSample(false, 1)</code>	Nondeterministic
<code>reduce(func)</code>	Combine the elements of the RDD together in parallel (e.g., sum).	<code>rdd.reduce((x, y) =&gt; x + y)</code>	9
<code>fold(zero)(func)</code>	Same as <code>reduce()</code> but with the provided zero value.	<code>rdd.fold(0)((x, y) =&gt; x + y)</code>	9



---

## KEY VALUE PAIRS

- ▶ `pairs = lines.map(lambda line: (line.split(" ")(0),x))`

---

# LOADING AND SAVING DATA

- ▶ Text Files :
  - ▶ `sc.textFile("file")` or `sc.wholeTextFiles("files")`
  - ▶ `sc.saveAsTextFile("outputFile")`