



Formação Engenheiro de Dados

RDD

RDD

- ◇ São recalculados a cada operação, não persistidos em memória (Lazy Evaluation)
 - ◇ Arquitetura pra grandes volumes de dados
- ◇ Para reuso, usar função RDD. persist() (diferentes formas, inclusive disco)

Lazy Evaluation



A Execução só começa quando é realmente necessária



Aplicado a transformações e carga de dados

Persistência

- ◊ Um RDD é imutável e recriado sempre que uma ação é executada
- ◊ Podemos persistir um RDD de diversas formas:
 - ◊ MEMORY_ONLY
 - ◊ MEMORY_ONLY_SER
 - ◊ MEMORY_AND_DISK
 - ◊ MEMORY_AND_DISK_SER
 - ◊ DISK_ONLY

Transformações e Ações

Transformações:

- Geram um novo RDD
- É mantida uma referência de dependências entre RDDs (lineage graph)

Ações:

- Cálculo exibido no console

Principais Transformações

Transformação	Descrição
<code>filter()</code>	Aplica uma filtro
<code>map()</code>	Aplica uma função
<code>sample()</code>	Gera subconjunto
<code>distinct()</code>	Retorna elementos únicos
<code>intersection()</code>	Retornar interseção de dois RDDs
<code>subtract()</code>	Subtrai conteúdo de um RDD
<code>cartesian()</code>	Gera o produto cartesiano de 2 RDDs
<code>union()</code>	Gera RDD com elementos de 2 RDDs

Principais Transformações para Chave- Valor

Transformação	Descrição
keys()	Retorna apenas as chaves
values()	Retorna apenas os valores
groupByKey()	Agrupar por chaves
sortByKey()	Ordena por chave

Principais Ações

Ação	Descrição
<code>collect()</code>	Retorna todo o RDD
<code>count()</code>	Conta o número de elementos
<code>countByValue()</code>	Contagem agrupada
<code>take()</code>	Retorna o número de elementos passados como parâmetro
<code>top()</code>	Retorna os primeiros elementos de acordo com o parâmetro
<code>reduce()</code>	Combina elementos usando computação paralela
<code>mean()</code>	Calcula a média
<code>sum()</code>	Calcula a soma
<code>min()</code>	Menor Valor
<code>max()</code>	Maior Valor
<code>variance()</code>	Calcula a Variância