



Data Science Academy

Como Funciona o HDFS?

Compreender a estrutura do HDFS é fundamental para o design de sistemas de processamento de dados altamente eficientes.

Vejamos alguns conceitos básicos do HDFS.

Blocos

O Hadoop “quebra” os arquivos recebidos em blocos e armazena-os de forma redundante em todo o cluster. Imagine um único arquivo grande, que é dividido em blocos, e os blocos são distribuídos entre os nodes disponíveis. Blocos HDFS são geralmente grandes, por padrão, 128 MB de tamanho, mas configuráveis pelo administrador do sistema.

WORM (Write Once Read Many times)

O Hadoop possui um paradigma diferente dos bancos de dados relacionais. Enquanto nos RDBMS os dados estão envolvidos em frequentes operações de leitura e escrita, no Hadoop o paradigma é de gravar um imenso conjunto de dados uma vez e realizar a leitura quantas vezes forem necessárias.

Replicação de dados

Os dados armazenados em todo o cluster são automaticamente replicados. Isso aumenta sua confiabilidade e disponibilidade. Por padrão, a replicação de arquivos é tríplice. O HDFS é otimizado para a leitura de grandes volumes de dados, em vez de leituras aleatórias.

Existem dois tipos de nós do cluster Hadoop:

NameNode: mantém o controle (metadados) de blocos que compõem um arquivo e também a localização desses blocos.

DataNodes: armazenam os blocos.

Vamos supor que temos três grandes arquivos: arquivo1, arquivo2, arquivo4, conforme ilustrado abaixo. Esses arquivos são divididos em blocos e espalhados por todos os nós do cluster. Por padrão, eles são armazenados de forma tríplice nos DataNodes. Informações sobre onde os



Data Science Academy

pedaços de arquivos podem ser encontrados são armazenados como metadados em um nó chamado NameNode.

NameNode armazena os metadados:

Metadado	<div>/home/hadoop/arquivo1 → 1, 2</div> <div>/home/hadoop/arquivo2 → 3, 4, 5</div> <div>/home/hadoop/arquivo3 → 6</div>
----------	---

DataNodes armazena os blocos de dados:

