

Resumo do que eu entendi.

O projeto mostra, de forma prática e didática, como configurar o ambiente Spark, manipular arquivos de log da NASA e realizar análises com PySpark — desde a leitura e filtragem de dados até a contagem de acessos, erros 404 e volume de tráfego — utilizando o poder do processamento distribuído.

Célula 1: Instala e configura o ambiente — atualiza pacotes, instala o Java (necessário para o Spark), instala o PySpark e clona um repositório do GitHub. Também mostra como manipular arquivos compactados (.gz), verificar tamanho, descompactar e visualizar conteúdo.

Célula 2: Explica o conceito de RDD (Resilient Distributed Dataset), a estrutura principal do Spark para processar grandes volumes de dados de forma distribuída, paralela, tolerante a falhas e imutável.

Célula 3: Mostra a configuração do SparkConf e SparkContext, definindo onde o Spark rodará, quanta memória usará e como o Python se conecta ao motor Spark.

Célula 4: Cria um RDD a partir de um arquivo de texto, usa cache para otimizar leituras e exibe amostras das linhas.

Células 5 a 6: Ensina a extrair o IP (primeiro token) das linhas e contar quantos hosts distintos acessaram o servidor nos logs de julho e agosto.

Célula 7: Define a função código404 para identificar linhas com erro HTTP 404 (página não encontrada), com tratamento de erros.

Célula 8: Conta quantos erros 404 ocorreram em julho e agosto, comparando métodos robustos e simples.

Célula 9: Extrai as URLs que causaram erro 404 e conta quantas vezes cada uma apareceu.

Célula 10: Ordena e exibe as 5 URLs com mais erros 404.

Célula 11: Analisa os erros 404 por dia, criando uma função para contar quantos ocorreram em cada data.

Célula 12: Ordena os resultados com lambda functions, exibindo os 10 dias com mais erros em ordem decrescente.

Célula 13: Calcula o total de bytes transferidos em julho e agosto, somando os valores dos logs.

Célula Final: Encerra o Spark e libera os recursos do sistema.