

Taller de Programación Apache Spark



Introducción

Apache Spark es un motor de procesamiento distribuido diseñado para manejar grandes volúmenes de datos de manera rápida y eficiente. Nació en el laboratorio AMP de la Universidad de California, Berkeley, como respuesta a las limitaciones de frameworks más antiguos como Hadoop MapReduce. Su principal fortaleza radica en su capacidad para realizar operaciones en memoria, lo que acelera significativamente el procesamiento de datos en comparación con los métodos tradicionales que dependen del almacenamiento en disco. Esta arquitectura hace que Spark sea ideal para tareas que requieren análisis complejos y rápidos, como el entrenamiento de modelos de machine learning, el procesamiento de flujos de datos en tiempo real y el análisis exploratorio a gran escala.

Con el tiempo, Spark ha evolucionado hasta convertirse en un ecosistema completo, integrando módulos para diversas necesidades. Spark SQL permite realizar consultas analíticas usando un enfoque declarativo similar al lenguaje SQL, facilitando su adopción por analistas y científicos de datos. Para el procesamiento de flujos en tiempo real, Spark Structured Streaming se ha convertido en un estándar, ofreciendo un enfoque unificado para trabajar tanto con datos estáticos como en movimiento. Además, su biblioteca MLlib proporciona herramientas avanzadas para machine learning, desde algoritmos clásicos como regresiones y árboles de decisión, hasta modelos más complejos que pueden escalar a grandes conjuntos de datos. Finalmente, GraphX amplía el potencial de Spark al permitir la creación y análisis de grafos, útiles para resolver problemas como la detección de comunidades o la búsqueda de caminos más cortos en redes complejas.

En 2025, Spark continúa siendo una referencia en el mundo del Big Data, con mejoras centradas en la eficiencia y la integración con tecnologías emergentes. Las últimas versiones han optimizado su motor Catalyst, lo que permite ejecutar consultas más rápidamente, y han potenciado su compatibilidad con el ecosistema de la nube, facilitando implementaciones escalables y seguras. Además, el soporte para lenguajes como Python mediante PySpark se ha fortalecido, integrando mejor con herramientas populares como Pandas y Polars, lo que hace que Spark sea aún más accesible para los científicos de datos. A su vez, la creciente adopción de tecnologías de inteligencia artificial generativa ha llevado a integrar más nativamente Spark con flujos de entrenamiento y despliegue de modelos, potenciando su rol en proyectos de data science.

Taller de Programación Apache Spark



Objetivos del Taller

El objetivo de este Taller es proporcionar a los profesionales de la ciencia de datos, ingeniería de datos y desarrollo de software una comprensión profunda y práctica de Apache Spark, permitiéndoles procesar y analizar grandes volúmenes de datos de manera eficiente y distribuida.

A lo largo del Taller, los participantes aprenderán a utilizar Spark para desarrollar aplicaciones de procesamiento de datos a gran escala, aprovechando su arquitectura optimizada para el procesamiento en memoria y su capacidad para ejecutar flujos de trabajo complejos de manera eficiente. Se abordarán temas clave como la creación y manipulación de DataFrames, el uso de Spark SQL para consultas avanzadas, el procesamiento en tiempo real mediante Spark Streaming y la aplicación de técnicas de machine learning con MLlib.

Este Taller está diseñado para que los participantes adquieran las habilidades necesarias para construir pipelines de datos eficientes y escalables, implementando buenas prácticas para el manejo de grandes conjuntos de datos. A través de ejercicios prácticos y casos reales, explorarán cómo optimizar el rendimiento de sus aplicaciones distribuidas mediante técnicas como la partición de datos, la caché en memoria y la paralelización de tareas.

Además, se pondrá énfasis en la integración de Apache Spark con otras herramientas del ecosistema Big Data, como Hadoop, Kafka y sistemas de almacenamiento en la nube, permitiendo a los participantes comprender su papel dentro de arquitecturas modernas de procesamiento de datos.

Al finalizar el Taller, los participantes estarán preparados para desarrollar aplicaciones de análisis de datos en tiempo real, implementar flujos de trabajo de ETL y construir modelos de machine learning distribuidos, dominando tanto los aspectos teóricos como prácticos del desarrollo con Apache Spark.

Taller de Programación Apache Spark

Contenidos



3

Introducción a Apache Spark

- ¿Qué es Apache Spark? Historia y evolución.
- Arquitectura de Spark: Driver, Executors y Cluster Manager.
- Diferencias entre Spark y Hadoop MapReduce.
- Casos de uso y aplicaciones en la industria.
- Instalación y configuración básica de Spark en Google Colaboratory.

Descargando e Instalando Spark en Google Colaboratory

- Configuración del entorno de trabajo.
- Instalación y configuración de Apache Spark en Colaboratory.
- Primer contacto con la CLI de Spark y ejecución de comandos básicos.
- Introducción a los entornos de ejecución: Local, YARN y Kubernetes.

Introducción a los RDD en Spark

- Concepto de RDD (Resilient Distributed Dataset).
- Creación de RDDs desde colecciones y archivos externos.
- Transformaciones y acciones en RDDs.
- Particionamiento y paralelismo en RDDs.
- Persistencia y almacenamiento de RDDs.

Transformaciones en un RDD

- Transformaciones básicas: `map()`, `filter()`, `flatMap()`.
- Transformaciones clave para el procesamiento distribuido: `groupBy()`, `reduceByKey()`, `aggregateByKey()`.
- Optimización del uso de transformaciones en grandes volúmenes de datos.
- Ejemplos prácticos de procesamiento con RDDs.

Acciones sobre un RDD en Spark

- Acciones básicas: `collect()`, `count()`, `first()`, `take()`.
- Acciones de reducción: `reduce()`, `fold()`, `aggregate()`.
- Acciones que implican escritura de datos: `saveAsTextFile()`, `saveAsSequenceFile()`.
- Comparación entre transformaciones y acciones para optimizar el rendimiento.

Aspectos Avanzados sobre RDD

- Uso de caché y persistencia para mejorar el rendimiento.
- Ajuste de particionamiento en RDDs para balancear carga de trabajo.
- Diagnóstico y solución de problemas de desempeño en RDDs.
- Análisis de ejecución con Spark UI.

Introducción a Spark SQL

- ¿Qué es Spark SQL y por qué utilizarlo?
- Creación y manipulación de DataFrames.
- Ejecución de consultas SQL sobre estructuras de datos distribuidas.
- Importación y exportación de datos en formatos CSV, JSON y Parquet.

Taller de Programación Apache Spark

Spark SQL Avanzado

- Optimización de consultas con el Catalyst Optimizer.
- Uso de UDFs (User Defined Functions) para extender las capacidades de Spark SQL.
- Creación de vistas temporales y gestión de tablas en memoria.
- Comparación entre Spark SQL y bases de datos tradicionales.

Funciones en Spark SQL

- Funciones de agregación: `sum()`, `avg()`, `count()`, `min()`, `max()`.
- Funciones de ventana (window functions) para análisis de series de tiempo.
- Funciones de transformación de datos: `regexp_replace()`, `split()`, `concat_ws()`.
- Creación de funciones definidas por el usuario (UDFs) en Python y Scala.

Conclusiones y Cierre

- Resumen de los conceptos aprendidos en el Taller.
- Preguntas y respuestas sobre los desafíos de Spark en entornos reales.
- Recomendaciones para seguir aprendiendo Spark y Big Data.



Taller de Programación Apache Spark

Metodología

El Taller de Apache Spark se desarrollará con un enfoque altamente práctico, combinando un 70% de ejercicios y laboratorios con un 30% de teoría. Esta estructura permitirá a los participantes no solo comprender los fundamentos de Spark, sino también aplicarlos en escenarios reales de procesamiento de datos a gran escala.

A lo largo del Taller, los participantes trabajarán en ejercicios prácticos y casos reales, utilizando Apache Spark para manipular y analizar grandes volúmenes de datos. Se enfocará en la aplicación de los conceptos teóricos mediante la implementación de scripts en PySpark y Scala, asegurando que los alumnos adquieran experiencia en la escritura y optimización de código distribuido. Cada módulo incluirá laboratorios diseñados para reforzar los conocimientos adquiridos, abordando desafíos comunes en el manejo de datos masivos.

El Taller fomentará la interacción y el aprendizaje colaborativo, incentivando a los alumnos a compartir sus experiencias, discutir soluciones y resolver problemas en equipo. La participación será clave para reforzar el conocimiento, permitiendo a los estudiantes enfrentar desafíos reales y resolverlos en un entorno guiado. Se trabajará con datasets reales y herramientas de monitorización de rendimiento, como Spark UI, para analizar la eficiencia del código ejecutado.

Además, el Taller incluirá proyectos integradores donde los participantes deberán aplicar todos los conocimientos adquiridos en el procesamiento de datos estructurados y no estructurados, optimización de consultas con Spark SQL, y desarrollo de flujos de trabajo de ETL escalables. Se brindará acceso a entornos de trabajo en la nube y local, asegurando que los alumnos experimenten con Spark en diferentes configuraciones.

La formación también incluirá evaluaciones continuas a través de pequeños retos de código y análisis de desempeño, garantizando que los estudiantes no solo comprendan los conceptos, sino que también sean capaces de aplicarlos eficientemente en un contexto profesional.

Al finalizar el Taller, los participantes estarán preparados para desarrollar soluciones de procesamiento de datos distribuidos utilizando Apache Spark en entornos empresariales, aplicando buenas prácticas para la optimización del rendimiento y la escalabilidad. La metodología está diseñada para proporcionar una experiencia de aprendizaje dinámica, práctica y orientada a la resolución de problemas reales en el ámbito del Big Data.

Taller de Programación Apache Spark



Requisitos

Para participar en este Taller de Apache Spark, es necesario contar con ciertos conocimientos y habilidades previas que facilitarán la comprensión de los conceptos y el desarrollo de las actividades prácticas.

Se espera que los participantes tengan una comprensión básica de programación, preferiblemente en Python o Scala, ya que serán los lenguajes principales utilizados en las implementaciones de Spark. También es importante contar con conocimientos fundamentales sobre conceptos de Big Data, como procesamiento distribuido y almacenamiento de datos a gran escala.

Además, se recomienda tener experiencia previa en el manejo de sistemas operativos tipo Unix/Linux, dado que muchas de las prácticas se realizarán en entornos de línea de comandos y mediante scripts ejecutados en terminal.

Familiaridad básica con bases de datos y SQL es esencial, ya que se trabajará intensivamente con Spark SQL para consultar y analizar grandes conjuntos de datos.

No es necesario tener experiencia previa en Apache Spark, pero es fundamental tener disposición para abordar retos prácticos y participar en actividades colaborativas que simularán situaciones reales en proyectos de procesamiento de datos a gran escala.

Este Taller está diseñado para que tanto profesionales que ya trabajan en proyectos de análisis de datos como aquellos que deseen adentrarse en el mundo del Big Data puedan adquirir las habilidades necesarias para utilizar Apache Spark de forma eficiente y productiva.

Dirigido a:

Profesionales de la Ciencia de Datos, Ingeniería de Datos y Desarrollo de Software

Generalidades

- Duración 32 horas cronológicas.
- Taller cerrado.