

Taller Fundamentos de Apache Livy



Introducción

En el ecosistema del Big Data y Apache Spark, uno de los principales desafíos es la ejecución y gestión remota de trabajos distribuidos. Apache Livy surge como una solución para facilitar la interacción con clústeres de Spark, permitiendo ejecutar consultas y trabajos sin necesidad de acceso directo a la infraestructura.

Apache Livy es un servicio REST para Apache Spark que permite enviar y administrar tareas de Spark desde aplicaciones externas, plataformas web o entornos como Jupyter Notebook y Apache Zeppelin. Su objetivo es proporcionar una interfaz sencilla y escalable para ejecutar código en Spark de forma remota, sin necesidad de usar herramientas tradicionales como spark-submit o interactuar directamente con los nodos del clúster.

Este servicio es especialmente útil en entornos multiusuario, donde varios clientes necesitan enviar consultas a un mismo clúster de Spark sin interferir en las ejecuciones de otros usuarios. Livy permite crear sesiones individuales, gestionar trabajos en segundo plano y recuperar resultados de manera eficiente, optimizando el uso de recursos en clústeres grandes.

Livy es compatible con Python (PySpark), Scala y Java, lo que lo convierte en una herramienta flexible para desarrolladores, científicos de datos y analistas que trabajan con grandes volúmenes de datos. Además, su integración con plataformas en la nube, como Azure Databricks, Amazon EMR y Cloudera, facilita su adopción en entornos empresariales donde Spark es la base del procesamiento de datos.

En este Taller básico, exploraremos los fundamentos de Apache Livy, incluyendo su instalación, configuración y uso para ejecutar tareas en un clúster de Spark. A lo largo del Taller, aprenderás a enviar consultas, gestionar sesiones y optimizar el rendimiento de tus trabajos, aprovechando las ventajas de este servicio REST.

Taller Fundamentos de Apache Livy



Objetivos del Taller

El objetivo de este Taller es proporcionar a los participantes una comprensión clara y práctica de Apache Livy, permitiéndoles utilizar este servicio REST para gestionar y ejecutar trabajos en Apache Spark de manera remota. A través de este Taller, los alumnos aprenderán a interactuar con clústeres de Spark sin necesidad de acceso directo a los nodos, optimizando el desarrollo y la administración de aplicaciones basadas en procesamiento distribuido.

El Taller está diseñado para introducir a los participantes en los conceptos esenciales de Livy, abarcando desde su instalación y configuración hasta la ejecución de consultas y la gestión de sesiones en entornos multiusuario. Se explorarán casos de uso en los que Livy simplifica la interacción con Spark en notebooks interactivos, aplicaciones web y entornos de ciencia de datos.

A lo largo del Taller, los alumnos adquirirán habilidades para enviar trabajos a Spark mediante la API REST de Livy, administrar sesiones concurrentes, recuperar resultados y optimizar la ejecución de consultas en entornos distribuidos. También se abordará su integración con plataformas como Jupyter Notebook, Apache Zeppelin y servicios en la nube como Azure Databricks y Amazon EMR.

Este Taller está dirigido a ingenieros de datos, científicos de datos, desarrolladores y administradores de sistemas que buscan mejorar la eficiencia en la ejecución de tareas en Spark. Al finalizar el Taller, los participantes estarán capacitados para usar Apache Livy en proyectos reales, facilitando la automatización y la administración de cargas de trabajo en clústeres de Apache Spark.

Taller Fundamentos de Apache Livy

Contenidos

Introducción a Apache Livy y su Arquitectura

- ¿Qué es Apache Livy y por qué usarlo?
- Diferencias entre ejecutar Spark localmente y con Livy.
- Arquitectura de Livy: cómo interactúa con Apache Spark.
- Casos de uso y aplicaciones en Big Data.

Instalación y Configuración de Apache Livy

- Requisitos previos para instalar Livy.
- Configuración de Livy en un entorno local y en la nube.
- Integración con Apache Spark y gestión de clústeres.
- Verificación de instalación y primeros comandos.

Uso de la API REST de Apache Livy

- Introducción a la API REST de Livy: solicitudes y respuestas en JSON.
- Creación y gestión de sesiones en Apache Livy.
- Envío de trabajos interactivos en PySpark y Scala.
- Recuperación de resultados y manejo de errores.

Integración de Livy con Notebooks y Aplicaciones

- Uso de Jupyter Notebook y Apache Zeppelin con Livy.
- Integración de Livy en aplicaciones web para análisis de datos.
- Configuración de seguridad y permisos en Livy.
- Consideraciones para el uso de Livy en entornos empresariales.

Proyecto Final y Mejores Prácticas

- Implementación de un flujo de trabajo con Apache Livy y Spark.
- Optimización de sesiones y administración de recursos en Livy.
- Monitoreo y resolución de problemas comunes.
- Conclusiones y próximos pasos en la adopción de Livy.

Taller Fundamentos de Apache Livy

Metodología

El Taller de Fundamentos de Apache Livy se desarrollará con un enfoque práctico y progresivo, combinando teoría con ejercicios aplicados en entornos reales de procesamiento distribuido. Aproximadamente el 60% del Taller estará basado en ejercicios prácticos, mientras que el 40% se enfocará en la teoría esencial para comprender la arquitectura y el uso de Livy con Apache Spark.

Cada sesión incluirá demostraciones en vivo y ejercicios guiados, donde los participantes aprenderán a instalar, configurar y utilizar Apache Livy para ejecutar trabajos en Spark de manera remota. Se trabajará con la API REST de Livy, permitiendo que los alumnos envíen tareas, administren sesiones y recuperen resultados en entornos distribuidos.

El Taller fomentará la participación, incentivando a los estudiantes a realizar pruebas en sus propios entornos y discutir posibles aplicaciones de Livy en distintos escenarios. Se proporcionarán ejemplos prácticos con Jupyter Notebook, Apache Zeppelin y herramientas de integración con Spark, asegurando que los alumnos adquieran habilidades aplicables en proyectos reales.

Para reforzar el aprendizaje, el Taller finalizará con un proyecto práctico, donde cada participante implementará un flujo de trabajo con Apache Livy y Spark, ejecutando consultas y analizando resultados mediante la API REST. Este enfoque permitirá a los alumnos consolidar sus conocimientos y desarrollar confianza en el uso de Livy para la gestión de trabajos en Spark.

Al finalizar el Taller, los participantes habrán adquirido una base sólida en Apache Livy, entendiendo su arquitectura, configurando entornos y ejecutando trabajos en Spark de manera remota. Esto les permitirá integrar Livy en sus proyectos de Big Data e Inteligencia Artificial, optimizando la ejecución de tareas en entornos empresariales y en la nube.

Taller Fundamentos de Apache Livy

Requisitos

Para aprovechar al máximo este Taller, se recomienda que los participantes tengan conocimientos básicos en Big Data, procesamiento distribuido y Apache Spark. Sin embargo, dado que es un Taller de nivel básico, se explicarán los conceptos esenciales antes de abordar la implementación práctica.

Es recomendable contar con experiencia en programación en Python, ya que estos serán los principales lenguajes utilizados para ejecutar consultas y trabajos en Apache Livy. Familiaridad con conceptos de procesamiento de datos y estructuras como DataFrames y RDDs en Spark será útil, aunque no obligatorio.

Se espera que los participantes tengan conocimientos básicos en el uso de línea de comandos en Linux, ya que Livy se configura y gestiona mediante comandos en entornos distribuidos. También se recomienda tener experiencia con APIs REST y JSON, ya que la interacción con Livy se realiza a través de solicitudes HTTP.

No es necesario tener experiencia previa en Apache Livy, pero sí es recomendable haber trabajado con Apache Spark en entornos locales o en la nube. Se sugiere contar con un entorno de desarrollo listo para pruebas, como un clúster de Spark local o en plataformas como Azure Databricks, Amazon EMR o Cloudera.

Este Taller está diseñado para ingenieros de datos, científicos de datos, desarrolladores y administradores de sistemas interesados en optimizar la ejecución de trabajos en Spark mediante una interfaz remota y escalable.

Dirigido a:

Ingenieros de Datos, Científicos de Datos, Desarrolladores y Administradores de Sistemas

Generalidades

- Duración 16 horas cronológicas.
- Taller cerrado.