

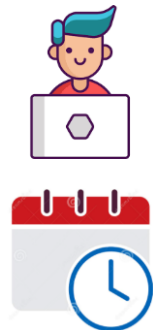


Bogotá R Users Group

BIG DATA CON R: SPARKLYR Y SPARKR

José Fernando Zea

febrero 11 de 2020



jfzeac@gmail.com

AGENDA

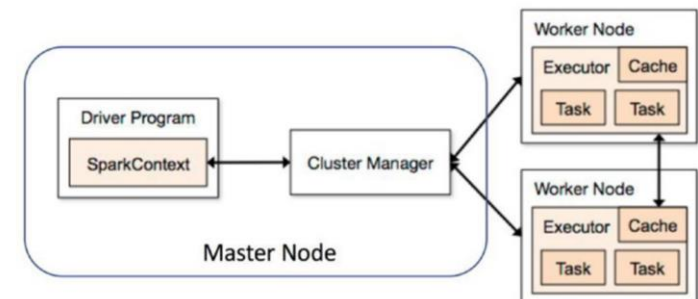
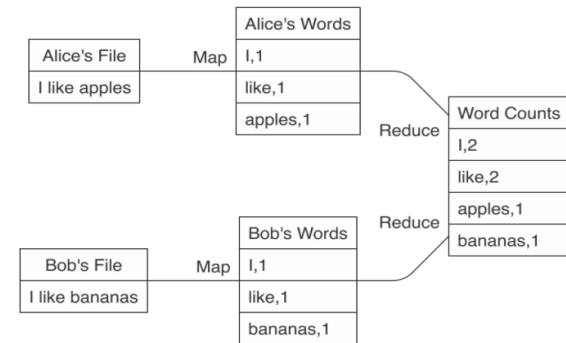
1. Contexto e historia
2. Puesta a punta de herramientas
3. Infraestructura
4. Configuración Databricks Community Edition
5. Ejemplos
6. ¿Cómo profundizo?

UN POCO DE HISTORIA ...

Hadoop y MapReduce (2004, 2006): Google, agregar y combinar. HDFS

Hive (2008): Facebook, MapReduce con SQL (Hive)

Spark (2009): Fundación Apache, datos en memoria.
Resilient Distributed Dataset



ANTES DEL WORKSHOP

1. Crear cuenta en databricks community edition:

<https://community.cloud.databricks.com/>

2. Crear *cluster* (colocar nombre a cluster. Usar configuración por defecto), esperar aprox 5 mins.

3. En el icono superior dar click y crear nuevo notebook: colocarle nombre y escoger de lenguaje R. El *cluster* creado aparece.

4. Importar Datos: ir al menú Data y a continuación picar en create table.

Verificar que se selecciona la pestaña Upload File y arrastar el archivo veredas.csv descargado de github (o dar click en browse y buscar el archivo).

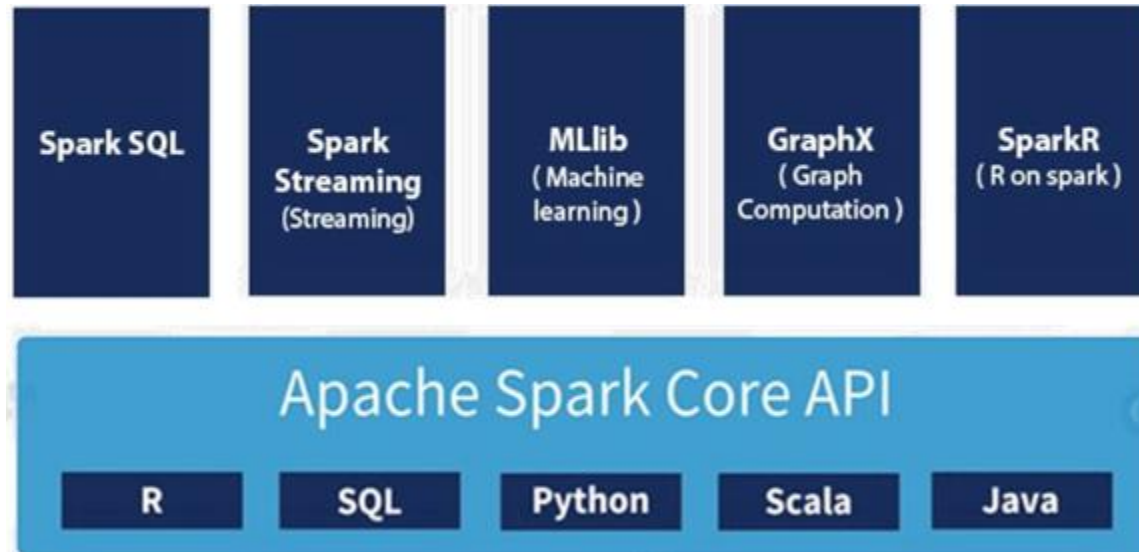
5. Vovler a nuestro notebook creado (En el menú Home)

0. Descargar de veredas.csv de:
<https://github.com/bogota-r/big-data-con-r>
<https://github.com/josezea/BogotaRUsers>

The image contains four numbered screenshots of the Databricks web interface, illustrating the setup process:

- 1:** The 'Clusters' page with the '+ Create Cluster' button highlighted by a blue arrow.
- 2:** The 'Clusters' page showing a cluster named 'rusers' in a 'Running' state.
- 3:** The 'Data' page with the 'Create Table' button highlighted by a blue arrow.
- 4:** The 'Data' page showing a message: 'You need to create a cluster to access tables'.

HERRAMIENTAS DE USO



SparkR



Fuente: Sing (2019), Lurashi et al (2019)

EJEMPLO

- Datos de las veredas del sector rural colombiano:

<https://geoportal.dane.gov.co/servicios/descarga-y-metadatos/descarga-nivel-de-referencia-de-veredas/>

- <https://github.com/bogota-r/big-data-con-r>

- <https://github.com/josezea/BogotaRUsers>

