

# Ciencia de Datos en Spark con *sparklyr* : : GUÍA RÁPIDA



## Conectar

### DATABRICKS CONNECT (v2)

1. Abre tu archivo .Renviron:
 

```
usethis::edit_r_environ()
```
2. En el archivo .Renviron agregue la dirección URL y el token del host de Databricks (PAT):
  - DATABRICKS\_HOST = [Tu dirección URL]
  - DATABRICKS\_TOKEN = [Tu token PAT]
3. Instalar extensión: `install.packages("pysparklyr")`
4. Abrir conexión:
 

```
sc <- spark_connect(
  cluster_id = "[Your cluster's ID]",
  method = "databricks_connect"
)
```

= Soportado en Databricks Connect v2

### CLÚSTER INDEPENDIENTE

1. Instale RStudio Server en uno de los nodos existentes o en un servidor de la misma LAN
2. Abrir una conexión
 

```
spark_connect(master="spark://host:port",
  version = "3.2",
  spark_home = [path to Spark])
```

### CLIENTE DE YARN

1. Instalación de RStudio Server en un nodo perimetral
2. Busque la ruta de acceso al directorio principal de Spark del clúster, normalmente es `"/usr/lib/spark"`
3. Ejemplo de configuración básica
 

```
conf <- spark_config()
conf$spark.executor.memory <- "300M"
conf$spark.executor.cores <- 2
conf$spark.executor.instances <- 3
conf$spark.dynamicAllocation.enabled <- "false"
```
4. Abrir una conexión
 

```
sc <- spark_connect(master = "yarn",
  spark_home = "/usr/lib/spark/",
  version = "2.1.0", config = conf)
```

### CLÚSTER DE YARN

1. Asegúrese de tener copias de los archivos `yarn-site.xml` y `hive-site.xml` en el servidor RStudio
2. Apunte las variables de entorno a las rutas correctas
 

```
Sys.setenv(JAVA_HOME="[Path]")
Sys.setenv(SPARK_HOME="[Path]")
Sys.setenv(YARN_CONF_DIR="[Path]")
```
3. Abrir una conexión
 

```
sc <- spark_connect(master = "yarn-cluster")
```

### KUBERNETES

1. Utilice lo siguiente para obtener el host y el puerto
 

```
system2("kubectl", "cluster-info")
```
2. Abrir conexión
 

```
sc <- spark_connect(config =
  spark_config_kubernetes(
    "k8s://https://[HOST]:[PORT]",
    account = "default",
    image = "docker.io/owner/repo:version"
  ))
```

### MODO LOCAL

No se requiere clúster. [Usar para solo para aprendizaje](#)

1. Instale una versión local de Spark: `spark_install()`
2. Abrir conexión
 

```
sc <- spark_connect(master="local")
```

### NUBE

Azure - `spark_connect(method = "synapse")`  
 Qubole - `spark_connect(method = "qubole")`

## Importar



### LEER UN ARCHIVO EN SPARK

Argumentos que se aplican a todas las funciones:  
`sc, name, path, options=list(), repartition=0, memory=TRUE, overwrite=TRUE`

<b>CSV</b>	<code>spark_read_csv(header = TRUE, columns=NULL, infer_schema=TRUE, delimiter = ",", quote = "\"", escape = "\\", charset = "UTF-8", null_value = NULL)</code>
<b>JSON</b>	<code>spark_read_json()</code>
<b>PARQUET</b>	<code>spark_read_parquet()</code>
<b>TEXT</b>	<code>spark_read_text()</code>
<b>DELTA</b>	<code>spark_read_delta()</code>

### DE UNA TABLA

`dplyr::tbl(scr, ...)` - Crea una referencia a la tabla sin cargar sus datos en la memoria  
`dbplyr::in_catalog()` - Habilita una dirección de tabla de tres partes  
`x <- tbl(sc, in_catalog("catalog", "schema", "table"))`

### Importar

- Desde R (`copy_to()`)
- Un archivo (`spark_read_*`)
- Una Hive table (`tbl()`)

### Visualizar

- Recopile el resultado, grafique en R

### Modelado

- Spark MLib (`m1_*`)
- Extensión H2O

### Manipular

- Verbos **dplyr**
- Comandos **tidyr**
- Transformador de características (`ft_*`)
- SQL directo de Spark (**DBI**)

[R for Data Science](#), [Wickham](#), [Çetinkaya-Rundel](#), [Grolemund](#)

### MARCO DE DATOS DE R EN SPARK

`dplyr::copy_to(dest, df, name)`

Apache Arrow acelera la transferencia de datos entre R y Spark. Para usarlo, simplemente cargue la biblioteca

```
library(sparklyr)
library(arrow)
```

## Manipular

### VERBOS DPLYR

Se traduce en instrucciones SQL de Spark

```
copy_to(sc, mtcars) |>
mutate(trm = ifelse(am == 0,
  "auto", "man")) |>
group_by(trm) |>
summarise_all(mean)
```

### TIDYR

- `pivot_longer()` - Contraer varias columnas en dos.
- `pivot_wider()` - Expanda dos columnas en varias.
- `nest()` / `unnest()` - Convierta grupos de celdas en columnas de lista y viceversa.
- `unite()` / `separate()` - Divida una sola columna en varias columnas y viceversa.
- `fill()` - Rellene NA con el valor anterior

### Communicar

Recolecte los resultados usando R y comunique con **Quarto**

### TRANSFORMADORES DE CARACTERÍSTICAS

- `ft_binarizer()` - Asigna valores basándose en un umbral
- `ft_bucketizer()` - De columna numérica a columna discretizada
- `ft_count_vectorizer()` - Extrae un vocabulario de un documento
- `ft_discrete_cosine_transform()` - Transformada discreta de coseno 1D de un vector real
- `ft_elementwise_product()` - Producto elemental entre 2 cols
- `ft_hashing_tf()` - Asigna una secuencia de términos a sus frecuencias de términos mediante el truco de hash.
- `ft_idf()` - Calcule la frecuencia inversa de documentos (IDF) dada una colección de documentos.
- `ft_imputer()` - El estimador de imputación para completar los valores faltantes, utiliza la media o la mediana de las columnas.
- `ft_index_to_string()` - Indexar etiquetas de nuevo para etiquetar como cadenas
- `ft_interaction()` - Toma las columnas Double y Vector y genera un vector aplanado de sus interacciones de entidades.
- `ft_max_abs_scaler()` - Cambie la escala de cada entidad individualmente al rango [-1, 1]
- `ft_min_max_scaler()` - Cambiar la escala de cada entidad a un rango común [mín., máx.] linealmente
- `ft_ngram()` - Convierte la matriz de cadenas de entrada en una matriz de n-gramas
- `ft_bucketed_random_projection_lsh()`  
`ft_minhash_lsh()` - Funciones hash sensibles a la localidad para la distancia euclidiana y la distancia de Jaccard (MinHash)

# Ciencia de Datos en Spark con *sparklyr* : : GUÍA RÁPIDA



- ft\_normalizer()** - Normalizar un vector para que tenga una norma unitaria usando la norma p dada
- ft\_one\_hot\_encoder()** - Continuo a vectores binarios
- ft\_pca()** - Proyecte vectores a un espacio dimensional inferior de los k componentes principales superiores.
- ft\_quantile\_discretizer()** - Continuo a valores categóricos agrupados.
- ft\_regex\_tokenizer()** - Extrae tokens mediante el patrón de expresiones regulares proporcionado para dividir el texto.
- ft\_robust\_scaler()** - Elimina la mediana y escala de acuerdo con la escala estándar.
- ft\_standard\_scaler()** - Elimina la media y escala a la varianza unitaria. Uso de estadísticas de resumen de columnas
- ft\_stop\_words\_remover()** - Filtra las palabras vacías de la entrada
- ft\_string\_indexer()** - Columna de etiquetas en una columna de índices de etiquetas.
- ft\_tokenizer()** - Convierte a minúsculas y luego lo divide por espacios en blanco
- ft\_vector\_assembler()** - Combinar vectores en un solo vector de fila
- ft\_vector\_indexer()** - Indexación de columnas de entidades categóricas en un conjunto de datos de Vector
- ft\_vector\_slicer()** - Toma un vector de características y genera un nuevo vector de características con una submatriz de las características originales
- ft\_word2vec()** - Word2Vec transforma una palabra en un código

## Modelado

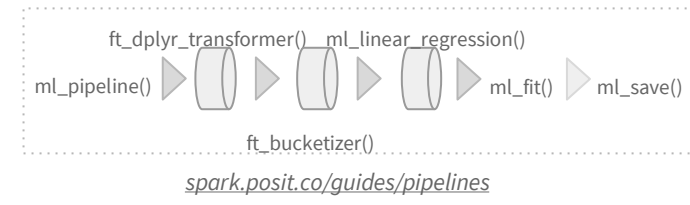
- REGRESIÓN**
  - ml\_linear\_regression()** - Regresión lineal.
  - ml\_aft\_survival\_regression()** - Modelo de regresión de supervivencia paramétrico denominado modelo de tiempo de falla acelerado (AFT)
  - ml\_generalized\_linear\_regression()** - GLM
  - ml\_isotonic\_regression()** - Utiliza el algoritmo de infractores adyacentes al grupo en paralelo.
  - ml\_random\_forest\_regressor()** - Regresión con bosques aleatorios.
- CLASIFICACIÓN**
  - ml\_linear\_svc()** - Clasificación mediante máquinas de vectores de soporte lineales
  - ml\_logistic\_regression()** - Regresión logística
  - ml\_multilayer\_perceptron\_classifier()** - Basado en el Perceptrón Multicapa.
  - ml\_naive\_bayes()** - Es compatible con Multinomial NB, que puede manejar datos discretos finitamente soportados
  - ml\_one\_vs\_rest()** - Reducción de Multiclase, realiza la reducción utilizando la estrategia de uno contra todos.
- ÁRBOL**
  - ml\_decision\_tree\_classifier()** | **ml\_decision\_tree()** | **ml\_decision\_tree\_regressor()** - Clasificación y regresión mediante árboles de decisión
  - ml\_gbt\_classifier()** | **ml\_gradient\_boosted\_trees()** | **ml\_gbt\_regressor()** - Clasificación binaria y regresión mediante árboles potenciados por gradiente
  - ml\_random\_forest\_classifier()** - Clasificación y regresión mediante bosques aleatorios.
  - ml\_feature\_importances()** | **ml\_tree\_feature\_importance()** - Importancia de las características para los modelos de árbol
- AGRUPAMIENTO**
  - ml\_bisecting\_kmeans()** - Un algoritmo de bisección de k-medias basado en el artículo
  - ml\_lda()** | **ml\_describe\_topics()** | **ml\_log\_likelihood()** | **ml\_log\_perplexity()** | **ml\_topics\_matrix()** - Modelo de tema LDA diseñado para documentos de texto.
  - ml\_gaussian\_mixture()** - Maximización de expectativas para modelos de mezclas gaussianas (GMM) multivariantes

- ml\_kmeans()** | **ml\_compute\_cost()** | **ml\_compute\_silhouette\_measure()** - Agrupación en clústeres con soporte para k-means
- ml\_power\_iteration()** - Para agrupar vértices de un grafo dadas similitudes por pares como propiedades de borde.
- RECOMENDACIÓN**
  - ml\_als()** | **ml\_recommend()** - Recomendación mediante la factorización de matrices de mínimos cuadrados alternos
- EVALUACIÓN**
  - ml\_clustering\_evaluator()** - Evaluador de clústeres
  - ml\_evaluate()** - Métricas de rendimiento de proceso
  - ml\_binary\_classification\_evaluator()** | **ml\_binary\_classification\_eval()** | **ml\_classification\_eval()** - Conjunto de funciones para calcular las métricas de rendimiento de los modelos de predicción.
- PATRÓN FRECUENTE**
  - ml\_fpgrowth()** | **ml\_association\_rules()** | **ml\_freq\_itemsets()** - Un algoritmo paralelo de crecimiento de FP para minar conjuntos de elementos frecuentes.
  - ml\_freq\_seq\_patterns()** | **ml\_prefixspan()** - Algoritmo PrefixSpan para minar conjuntos de elementos frecuentes.
- ESTADÍSTICAS**
  - ml\_summary()** - Extrae una métrica del objeto de resumen de un modo de Spark MLL
  - ml\_corr()** - Calcular matriz de correlación
- CARACTERÍSTICA**
  - ml\_chisquare\_test(x, features, label)** - Prueba de independencia de Pearson para característica-etiqueta
  - ml\_default\_stop\_words()** - Carga las palabras vacías predeterminadas para el idioma especificado
- UTILIDADES**
  - ml\_call\_constructor()** - Identifica el constructor de ML de Sparklyr asociado para la JVM
  - ml\_model\_data()** - Extrae datos asociados a un modelo de Spark ML
  - ml\_standardize\_formula()** - Genera una cadena de fórmula a partir de las entradas del usuario
  - ml\_uid()** - Extrae el UID de un objeto de ML.

## Canalizaciones de ML

Cree fácilmente modelos formales de Spark Pipeline con R. Guarde la canalización en Sacala nativo. No tendrá dependencias de R.

- INICIALIZAR Y ENTRENAR**
  - ml\_pipeline()** - Inicializa una nueva canalización de Spark
  - ml\_fit()** - Entrena el modelo y genera un modelo de canalización de Spark.
- GUARDAR Y RECUPERAR**
  - ml\_save()** - Guarda en un formato que puede ser leído por Scala y PySpark.
  - ml\_read()** - Lee el objeto Spark en sparklyr.



## R distribuido

Ejecute código R arbitrario a escala dentro del clúster con **spark\_apply()**. Útil cuando se necesita una funcionalidad que solo está disponible en R, y para resolver "problemas vergonzosamente paralelos"

**spark\_apply(x, f, columns = NULL, memory = TRUE, group\_by = NULL, name = NULL, barrier = NULL, fetch\_result\_as\_sdf = TRUE)**

```
copy_to(sc, mtcars) |>
spark_apply(
  nrow, # R only function
  group_by = "am",
  columns = "am double, x long"
)
```

### Más información

spark.posit.co      therinspark.com

## Visualizar

### DPLYR + GGLOT2

```
group_by(cyl) |>
summarise(mpg_m = mean(mpg)) |>
collect() |>
ggplot() +
geom_col(aes(cyl, mpg_m))
```

Resumir en Spark  
Recopilar resultados en R  
Crear gráfica

