

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



## Intro

**sparklyr** es una interfaz de R para **Apache Spark™**. **sparklyr** nos permite escribir todo el código de nuestro análisis en R, pero con el verdadero procesamiento real sucediendo dentro de los clusters de Spark. Manipula y modela a gran escala fácilmente usando R y Spark mediante **sparklyr**.

## Importar



### LEE UN ARCHIVO EN SPARK

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

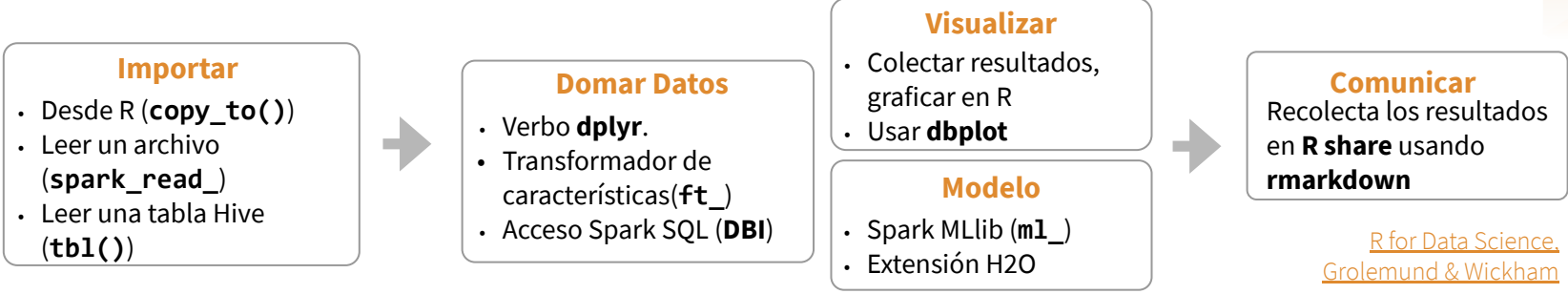
CSV	spark_read_csv( header = TRUE, columns=NULL, infer_schema=TRUE, delimiter = ",", quote = "\"", escape = "\\\"", charset = "UTF-8", null_value = NULL)
JSON	spark_read_json()
PARQUET	spark_read_parquet()
TEXTO	spark_read_text()
TABLA HIVE	spark_read_table()
ORC	spark_read_orc()
LIBSVM	spark_read_libsvm()
JDBC	spark_read_jdbc()
DELTA	spark_read_delta()

### COPIAR UN DATA FRAME DE R A SPARK

dplyr::copy\_to(dest, df, name)

### DE UNA TABLA EN HIVE

dplyr::tbl(scr, ...) Crea una referencia a la tabla sin cargarla a la memoria



## Domar Datos

### VERBOS DPLYR

Traduce a sentencias Spark SQL

copy\_to(sc, mtcars) %>%  
mutate(trm = ifelse(am == 0, "auto", "man")) %>%  
group\_by(trm) %>%  
summarise\_all(mean)

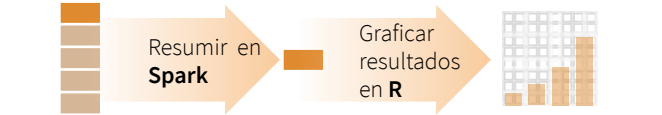
### TRANSFORMADORES DE CARACTERÍSTICAS

	ft_binarizer() - Valores asignados basados en el límite
	ft_bucketizer() - Columna numérica a columna discretizada
	ft_count_vectorizer() - Extrae vocabulario del documento
	ft_discrete_cosine_transform() - Coseno discreta a vector real
	ft_elementwise_product() - Producto basado en elementos entre 2 columnas
	ft_hashing_tf() - Asigna una secuencia de términos a sus frecuencias de términos utilizando el truco de hashing
	ft_idf() - Calcular la frecuencia de documentos inversa (IDF) dada una colección de documentos
	ft_imputer() - Estimador de imputación para completar valores perdidos, utiliza la media o la mediana de las columnas
	ft_index_to_string() - Index labels back to label as strings
	ft_interaction() - Toma columnas del tipo Doble y Vector y devuelve un vector achatado de sus interacciones características

	ft_max_abs_scaler() - Reescalar cada función individualmente al rango [-1, 1]
	ft_min_max_scaler() - Reescalar la escala de cada función individualmente a un rango común [min, max] linealmente
	ft_ngram() - Convierte la matriz de entrada de strings en una matriz de n-gramas.
	ft_bucketed_random_projection_lsh() - Funciones de Hashing sensibles a la localidad para distancia euclidiana y distancia Jaccard (MinHash)
	ft_normalizer() - Normaliza un vector para que tenga norma de unidad usando la p-norm dada
	ft_one_hot_encoder() - Vectores continuos a binarios
	ft_pca() - Proyecta vectores a un espacio dimensional menor de los top k componentes principales
	ft_quantile_discretizer() - Valores categóricos continuos a agrupados
	ft_regex_tokenizer() - Extrae tokens usando el patrón regex provisto para dividir el texto
	ft_standard_scaler() - Elimina la media y la escala a la varianza de la unidad utilizando estadísticas de resumen de columna
	ft_stop_words_remover() - Filtra las palabras vacías del input
	ft_string_indexer() - Columna de etiquetas en una columna de índices de etiquetas
	ft_tokenizer() - Convierte a minúsculas y luego los separa por espacios en blanco

	ft_vector_assembler() - Combina vectores en una sola fila-vector
	ft_vector_indexer() - Indexa columnas de caract. categóricas en un dataset de Vector
	ft_vector_slicer() - Toma un vector de características y genera uno nuevo con un subconjunto de las caract. originales
	ft_word2vec() - Word2Vec transforma una palabra en código

## Visualizar



### DPLYR + GGLOT2

copy\_to(sc, mtcars) %>%  
group\_by(cyl) %>%  
summarise(mpg\_m = mean(mpg)) %>%  
collect() %>%  
ggplot() +  
geom\_col(aes(cyl, mpg\_m))

Resumir en Spark  
Recolectar resultados en R  
Crear gráfico

### DBPLOT

copy\_to(sc, mtcars) %>%  
dbplot\_histogram(mpg) +  
labs(title = "Histogram of MPG")

dbplot\_histogram(data, x, bins = 30, binwidth = NULL) - Calcula los bins del histograma en Spark y grafica en ggplot2  
dbplot\_raster(data, x, y, fill = n(), resolution = 100, complete = FALSE) - Visualiza 2 variables continuas. Usar en lugar de geom\_point()

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



## Modelado

### REGRESIÓN

**ml\_linear\_regression()** - Regresión lineal  
**ml\_aft\_survival\_regression()** - Modelo paramétrico de regresión de supervivencia llamado AFT (por las siglas en inglés *Accelerated Failure Time*)  
**ml\_generalized\_linear\_regression()** - Modelo generalizado de regresión lineal  
**ml\_isotonic\_regression()** - Actualmente implementados usando el algoritmo PAVA. Sólo soporta algoritmos univariados (única característica)  
**ml\_random\_forest\_regressor()** - Regresión usando *Random Forests* (bosques aleatorios).

### CLASIFICACIÓN

**ml\_linear\_svc()** - Clasificación usando máquinas de vectores de soporte lineales  
**ml\_logistic\_regression()** - regresión logística  
**ml\_multilayer\_perceptron\_classifier()** - Modelo de clasificación basado en el Perceptrón Multicapa  
**ml\_naive\_bayes()** - Clasificadores Bayesianos Naive. Es compatible con Multinomial NB que puede manejar datos discretos con soporte finito  
**ml\_one\_vs\_rest()** - Reducción de clasificación multiclase a clasificación binaria, usando la estrategia “uno contra todos”.

### ÁRBOLES

**ml\_decision\_tree\_classifier()** | **ml\_decision\_tree()** | **ml\_decision\_tree\_regressor()** - clasificación y regresión utilizando árboles de decisión  
**ml\_gbt\_classifier()** | **ml\_gradient\_boosted\_trees()** | **ml\_gbt\_regressor()** - Clasificación binaria y regresión usando *Gradient Boosted Trees* (árboles de potenciación del gradiente)  
**ml\_random\_forest\_classifier()** - Clasificación y regresión usando *Random Forests*.  
**ml\_feature\_importances(model,...)** | **ml\_tree\_feature\_importance(model)** - Importancia de las características para los modelos de árboles

### AGRUPAMIENTO (*Clustering*)

**ml\_bisecting\_kmeans()** - Algoritmo *bisecting K-means* (bisección de k-medias) basado en la publicación.  
**ml\_lda()** | **ml\_describe\_topics()** | **ml\_log\_likelihood()** | **ml\_log\_perplexity()** | **ml\_topics\_matrix()** - LDA topic model diseñada para documentos de texto  
**ml\_gaussian\_mixture()** - Maximización de expectativas para los Modelos Gaussianos Multivariados (GMM)  
**ml\_kmeans()** | **ml\_compute\_cost()** - Agrupamiento de k-medias con soporte para k-medias

### CRECIMIENTO FP

**ml\_fpgrowth()** | **ml\_association\_rules()** | **ml\_freq\_itemsets()** - Algoritmo paralelo de crecimiento FP para extraer conjuntos de elementos frecuentes.

### CARACTERÍSTICAS (*Features*)

**ml\_chisquare\_test(x,features,label)** - Prueba de independencia de Pearson para cada característica contra la etiqueta  
**ml\_default\_stop\_words()** - Carga las palabras vacías predeterminadas para el idioma dado

### STATS

**ml\_summary()** - Extrae una métrica desde el objeto resumen de un modelo Spark ML  
**ml\_corr()** - Calcula matriz de correlación

El paquete **correlate** se integra con **sparklyr**



```
copy_to(sc, mtcars) %>%  
correlate() %>%  
rplot()
```



### RECOMENDACIÓN

**ml\_als()** | **ml\_recommend()** - Recomendación usando la factorización de la matriz de Cuadrados Mínimos Alternos (*ALS- Alternating Least Squares*)

### EVALUACIÓN

**ml\_clustering\_evaluator()** - Evaluador para *clustering*  
**ml\_evaluate()** - Calcular métricas de rendimiento  
**ml\_binary\_classification\_evaluator()** | **ml\_binary\_classification\_eval()** | **ml\_classification\_eval()** - Un conjunto de funciones para calcular métricas de rendimiento para modelos de predicción.

### UTILIDADES

**ml\_call\_constructor()** - Identifica el constructor ML de **sparklyr** asociado para la JVM (Máquina Virtual Java)  
**ml\_model\_data()** - Extrae datos asociados a un modelo Spark ML

**ml\_standardize\_formula()** - Genera un string de fórmula usando *inputs* de usuarios, para ser usado en el constructor “ **ml\_model** ”  
**ml\_uid()** - Extrae el UID de un objeto ML

## Comenzar una sesión de Spark

### YARN CLIENT

1. Instale RStudio Server en uno de los nodos existentes, preferiblemente un nodo límite.
2. Localice la ruta al directorio principal de Spark del cluster, 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. Abra una conexión (algunas configuraciones básicas incluidas en el ejemplo)  
`sc <- spark_connect(master = "yarn",  
spark_home = "/usr/lib/spark/",  
version = "2.1.0", config = conf)`

### YARN CLUSTER

1. Asegúrese de tener copias de los archivos `yarn-site.xml` y `hive-site.xml` en el servidor RStudio
2. Señale variables de entorno a las rutas correctas  
`Sys.setenv(JAVA_HOME="[Path]")`  
`Sys.setenv(SPARK_HOME="[Path]")`  
`Sys.setenv(YARN_CONF_DIR="[Path]")`
3. Abra una conexión  
`sc <- spark_connect(master = "yarn-cluster")`

### CLUSTER INDEPENDIENTE

1. Instale RStudio Server en uno de los nodos existentes o en un servidor en la misma LAN
2. Instale una versión local de Spark:  
`spark_install (version = "2.0.1")`
3. Abra una conexión:  
`spark_connect(master="spark://host:port",  
version = "2.0.1",  
spark_home = spark_home_dir())`

### MODO LOCAL

No se requiere cluster. Usar sólo con fines de aprendizaje

1. Instale una versión local de Spark  
`spark_install("2.3")`
2. Abra una conexión  
`sc <- spark_connect(master="local")`

### KUBERNETES

1. Use lo siguiente para obtener el host y el puerto  
`system2("kubectl", "cluster-info")`
2. Abra una conexión  
`sc <- spark_connect(config =  
spark_config_kubernetes(  
"k8s://https://[HOST]>:[PORT]",  
account = "default",  
image = "docker.io/owner/repo:version",  
version = "2.3.1"))`

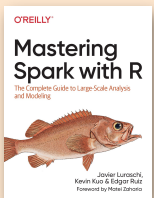
### MESOS

1. Instale Rstudio Server en uno de los nodos
2. Abra una conexión  
`sc <- spark_connect(master="[Mesos URL]")`

### CLOUD

**Databricks** - `spark_connect(method = "databricks")`  
**Qubole** - `spark_connect(method = "qubole")`

## Más información



[spark.rstudio.com](http://spark.rstudio.com)

[therinspark.com](http://therinspark.com)

