

# Data Science in Spark with *sparklyr* : : CHEAT SHEET



## Connect

### DATABRICKS CONNECT (v2)

1. Open your .Renviron file: `usethis::edit_r_environ()`
2. In the .Renviron file add your Databricks Host Url and Token (PAT):
  - `DATABRICKS_HOST = [Your Host URL]`
  - `DATABRICKS_TOKEN = [Your PAT]`
3. Install extension: `install.packages("pysparklyr")`
4. Open connection:
 

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

= Supported in Databricks Connect v2

### STANDALONE CLUSTER

1. Install RStudio Server on one of the existing nodes or a server in the same LAN
2. Open a connection
 

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

### YARN CLIENT

1. Install RStudio Server on an edge node
2. Locate path to the cluster's Spark Home Directory, it normally is `"/usr/lib/spark"`
3. Basic configuration example
 

```
conf <- spark_config()
conf$spark.executor.memory <- "300M"
conf$spark.executor.cores <- 2
conf$spark.executor.instances <- 3
conf$spark.dynamicAllocation.enabled <- "false"
```
4. Open a connection
 

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

### YARN CLUSTER

1. Make sure to have copies of the `yarn-site.xml` and `hive-site.xml` files in the RStudio Server
2. Point environment variables to the correct paths
 

```
Sys.setenv(JAVA_HOME="[Path]")
Sys.setenv(SPARK_HOME="[Path]")
Sys.setenv(YARN_CONF_DIR="[Path]")
```
3. Open a connection
 

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

### KUBERNETES

1. Use the following to obtain the Host and Port
2. Open a connection
 

```
system2("kubectl", "cluster-info")
sc <- spark_connect(config =
  spark_config_kubernetes(
    "k8s://https://[HOST]:[PORT]",
    account = "default",
    image = "docker.io/owner/repo:version"
  ))
```

### LOCAL MODE

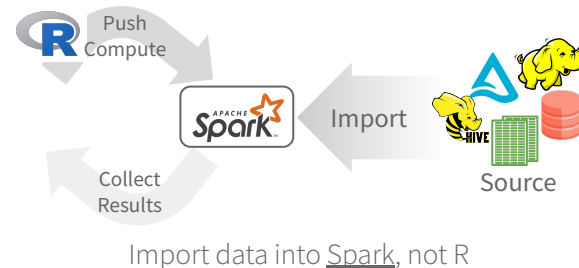
- No cluster required. [Use for learning purposes only](#)
1. Install a local version of Spark: `spark_install()`
  2. Open a connection
 

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

### CLOUD

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

## Import



### READ A FILE INTO SPARK

Arguments that apply to all functions:  
`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>

### FROM A TABLE

`dplyr::tbl(sc, ...)` - Creates a reference to the table without loading its data into memory  
`dbplyr::in_catalog()` - Enables a three part table address  
`x <- tbl(sc, in_catalog("catalog", "schema", "table"))`

### Import

- From R (`copy_to()`)
- Read a file (`spark_read_*`)
- Read Hive table (`tbl()`)

### Wrangle

- `dplyr` verb
- `tidyr` commands
- Feature transformer (`ft_*`)
- Direct Spark SQL (**DBI**)

[R for Data Science, Grolemund & Wickham](#)

### Visualize

- Collect result, plot in R

### Model

- Spark MLlib (`m1_*`)
- H2O Extension

### Communicate

Collect results into R share using **Quarto**

### R DATA FRAME INTO SPARK

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

Apache Arrow accelerates data transfer between R and Spark. To use, simply load the library

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

## Wrangle

### DPLYR VERBS

Translates into Spark SQL statements

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

### TIDYR

- `pivot_longer()` - Collapse several columns into two.
- `pivot_wider()` - Expand two columns into several.
- `nest() / unnest()` - Convert groups of cells into list-columns, and vice versa.
- `unite() / separate()` - Split a single column into several columns, and vice versa.
- `fill()` - Fill NA with the previous value

### FEATURE TRANSFORMERS

- `ft_binarizer()` - Assigned values based on threshold
- `ft_bucketizer()` - Numeric column to discretized column
- `ft_count_vectorizer()` - Extracts a vocabulary from document
- `ft_discrete_cosine_transform()` - 1D discrete cosine transform of a real vector
- `ft_elementwise_product()` - Element-wise product between 2 cols
- `ft_hashing_tf()` - Maps a sequence of terms to their term frequencies using the hashing trick.
- `ft_idf()` - Compute the Inverse Document Frequency (IDF) given a collection of documents.
- `ft_imputer()` - Imputation estimator for completing missing values, uses the mean or the median of the columns.
- `ft_index_to_string()` - Index labels back to label as strings
- `ft_interaction()` - Takes in Double and Vector columns and outputs a flattened vector of their feature interactions.
- `ft_max_abs_scaler()` - Rescale each feature individually to range [-1, 1]
- `ft_min_max_scaler()` - Rescale each feature to a common range [min, max] linearly
- `ft_ngram()` - Converts the input array of strings into an array of n-grams
- `ft_bucketed_random_projection_lsh()`  
`ft_minhash_lsh()` - Locality Sensitive Hashing functions for Euclidean distance and Jaccard distance (MinHash)

# Data Science in Spark with *sparklyr* : : CHEAT SHEET



- ft\_normalizer()** - Normalize a vector to have unit norm using the given p-norm
- ft\_one\_hot\_encoder()** - Continuous to binary vectors
- ft\_pca()** - Project vectors to a lower dimensional space of top k principal components.
- ft\_quantile\_discretizer()** - Continuous to binned categorical values.
- ft\_regex\_tokenizer()** - Extracts tokens either by using the provided regex pattern to split the text.
- ft\_robust\_scaler()** - Removes the median and scales according to standard scale.
- ft\_standard\_scaler()** - Removes the mean and scaling to unit variance using column summary statistics
- ft\_stop\_words\_remover()** - Filters out stop words from input
- ft\_string\_indexer()** - Column of labels into a column of label indices.
- ft\_tokenizer()** - Converts to lowercase and then splits it by white spaces
- ft\_vector\_assembler()** - Combine vectors into single row-vector
- ft\_vector\_indexer()** - Indexing categorical feature columns in a dataset of Vector
- ft\_vector\_slicer()** - Takes a feature vector and outputs a new feature vector with a subarray of the original features
- ft\_word2vec()** - Word2Vec transforms a word into a code

## Modeling

### REGRESSION

- ml\_linear\_regression()** - Linear regression.
- ml\_aft\_survival\_regression()** - Parametric survival regression model named accelerated failure time (AFT) model
- ml\_generalized\_linear\_regression()** - GLM
- ml\_isotonic\_regression()** - Uses parallelized pool adjacent violators algorithm.
- ml\_random\_forest\_regressor()** - Regression using random forests.

### CLASSIFICATION

- ml\_linear\_svc()** - Classification using linear support vector machines
- ml\_logistic\_regression()** - Logistic regression
- ml\_multilayer\_perceptron\_classifier()** - Based on the Multilayer Perceptron.
- ml\_naive\_bayes()** - It supports Multinomial NB which can handle finitely supported discrete data
- ml\_one\_vs\_rest()** - Reduction of Multiclass, performs reduction using one against all strategy.

### TREE

- ml\_decision\_tree\_classifier()** | **ml\_decision\_tree()** | **ml\_decision\_tree\_regressor()** - Classification and regression using decision trees
- ml\_gbt\_classifier()** | **ml\_gradient\_boosted\_trees()** | **ml\_gbt\_regressor()** - Binary classification and regression using gradient boosted trees
- ml\_random\_forest\_classifier()** - Classification and regression using random forests.
- ml\_feature\_importances()** | **ml\_tree\_feature\_importance()** - Feature Importance for Tree Models

### CLUSTERING

- ml\_bisecting\_kmeans()** - A bisecting k-means algorithm based on the paper
- ml\_lda()** | **ml\_describe\_topics()** | **ml\_log\_likelihood()** | **ml\_log\_perplexity()** | **ml\_topics\_matrix()** - LDA topic model designed for text documents.
- ml\_gaussian\_mixture()** - Expectation maximization for multivariate Gaussian Mixture Models (GMMs)

- ml\_kmeans()** | **ml\_compute\_cost()** | **ml\_compute\_silhouette\_measure()** - Clustering with support for k-means
- ml\_power\_iteration()** - For clustering vertices of a graph given pairwise similarities as edge properties.

### RECOMMENDATION

- ml\_als()** | **ml\_recommend()** - Recommendation using Alternating Least Squares matrix factorization

### EVALUATION

- ml\_clustering\_evaluator()** - Evaluator for clustering
- ml\_evaluate()** - Compute performance metrics
- ml\_binary\_classification\_evaluator()** | **ml\_binary\_classification\_eval()** | **ml\_classification\_eval()** - A set of functions to calculate performance metrics for prediction models.

### FREQUENT PATTERN

- ml\_fpgrowth()** | **ml\_association\_rules()** | **ml\_freq\_itemsets()** - A parallel FP-growth algorithm to mine frequent itemsets.
- ml\_freq\_seq\_patterns()** | **ml\_prefixspan()** - PrefixSpan algorithm for mining frequent itemsets.

### STATS

- ml\_summary()** - Extracts a metric from the summary object of a Spark ML model
- ml\_corr()** - Compute correlation matrix

### RECOMMENDATION

- ml\_als()** | **ml\_recommend()** - Recommendation using Alternating Least Squares matrix factorization

### FEATURE

- ml\_chisquare\_test(x, features, label)** - Pearson's independence test for every feature against the label
- ml\_default\_stop\_words()** - Loads the default stop words for the given language

### UTILITIES

- ml\_call\_constructor()** - Identifies the associated sparklyr ML constructor for the JVM
- ml\_model\_data()** - Extracts data associated with a Spark ML model
- ml\_standardize\_formula()** - Generates a formula string from user inputs
- ml\_uid()** - Extracts the UID of an ML object.

## ML Pipelines

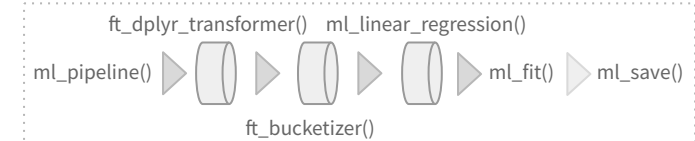
Easily create a formal Spark Pipeline models using R. Save the Pipeline in native Sacala. It will have **no dependencies on R**.

### INITIALIZE AND TRAIN

- ml\_pipeline()** - Initializes a new Spark Pipeline
- ml\_fit()** - Trains the model, outputs a Spark Pipeline Model.

### SAVE AND RETRIEVE

- ml\_save()** - Saves into a format that can be read by Scala and PySpark .
- ml\_read()** - Reads Spark object into sparklyr.



[spark.posit.co/guides/pipelines](https://spark.posit.co/guides/pipelines)

## Distributed R

Run arbitrary R code at scale inside your cluster with **spark\_apply()**. Useful when there you need functionality only available in R, and to solve 'embarrassingly parallel problems'

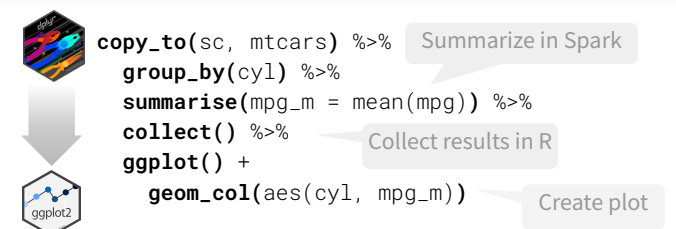
**spark\_apply(x, f, columns = NULL, memory = TRUE, group\_by = NULL, name = NULL, barrier = NULL, fetch\_result\_as\_sdf = TRUE)**

```

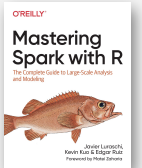
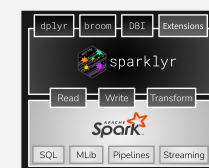
copy_to(sc, mtcars) %>%
  spark_apply(
    nrow, # R only function
    group_by = "am",
    columns = "am double, x long"
  )
  
```

## Visualize

### DPLYR + GGLOT2



### More Info



[spark.posit.co](https://spark.posit.co)

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