

Reactivity

COPY A DATA FRAME INTO SPARK

```
sdf_copy_to(sc, iris, "spark_iris")
```

```
sdf_copy_to(sc, x, name, memory, repartition,
overwrite)
```

IMPORT INTO SPARK FROM A FILE

Arguments that apply to all functions:
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()**

SPARK SQL COMMANDS

```
DBI::dbWriteTable(sc, "spark_iris", iris)
```

```
DBI::dbWriteTable(conn, name,
value)
```

FROM A TABLE IN HIVE

```
my_var <- tbl_cache(sc, name =
"hive_iris")
```

tbl_cache() (sc, name, force = TRUE)
Loads the table into memory

```
my_var <- dplyr::tbl(sc,
name = "hive_iris")
dplyr::tbl(sc, ...)
```

Creates a reference to the table
without loading it into memory

Wrangle

SPARK SQL VIA DPLYR VERBS

Translates into Spark SQL statements

```
my_table <- my_var %>%
filter(Species=="setosa") %>%
sample_n(10)
```

DIRECT SPARK SQL COMMANDS

```
my_table <- DBI::dbGetQuery(sc, "SELECT *
FROM iris LIMIT 10")
```

```
DBI::dbGetQuery(conn, statement)
```

SCALA API VIA SDF FUNCTIONS

```
sdf_mutate(.data)
```

Works like dplyr mutate function

```
sdf_partition(x, ..., weights = NULL, seed =
sample(.Machine$integer.max, 1))
sdf_partition(x, training = 0.5, test = 0.5)
```

```
sdf_register(x, name = NULL)
```

Gives a Spark DataFrame a table name

```
sdf_sample(x, fraction = 1, replacement =
TRUE, seed = NULL)
```

```
sdf_sort(x, columns)
```

Sorts by >= 1 columns in ascending order

```
sdf_with_unique_id(x, id = "id")
```

```
sdf_predict(object, newdata)
```

Spark DataFrame with predicted values

Visualize & Communicate

DOWNLOAD DATA TO R MEMORY

```
r_table <- collect(my_table)
plot(Petal.Width~Petal.Length, data=r_table)
dplyr::collect(x)
```

Download a Spark DataFrame to an R DataFrame

```
sdf_read_column(x, column)
```

Returns contents of a single column to R

SAVE FROM SPARK TO FILE SYSTEM

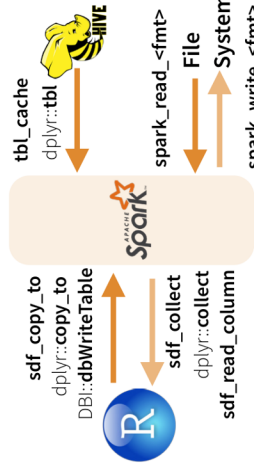
Arguments that apply to all functions: **x**, **path**

CSV **spark_read_csv()** (header = TRUE,
delimiter = ",", quote = "\"", escape = "\\",
charset = "UTF-8", null_value = NULL)

JSON **spark_read_json()** (mode = NULL)

PARQUET **spark_read_parquet()** (mode = NULL)

Reading & Writing from Apache Spark



Extensions

Create an R package that calls the full Spark API & provide interfaces to Spark packages.

CORE TYPES

spark_connection() Connection between R and the Spark shell process

spark_job() Instance of a remote Spark object

spark_dataframe() Instance of a remote Spark DataFrame object

CALL SPARK FROM R

invoke() Call a method on a Java object

invoke_new() Create a new object by invoking a constructor

invoke_static() Call a static method on an object

MACHINE LEARNING EXTENSIONS

ml_create_dummy_variables() **ml_options()**

ml_prepare_dataframe() **ml_model()**

ml_prepare_response_features_intercept()

Model (MLlib)

```
ml_decision_tree(my_table,
response = "Species", features =
c("Petal.Length", "Petal.Width"))
```

ml_als_factorization() (x, user.column = "user",
rating.column = "rating", item.column = "item",
rank = 10L, regularization.parameter = 0.1, iter.max = 10L,
ml.options = ml_options())

ml_decision_tree() (x, response, features, max.bins = 32L, max.depth = 5L, type = c("auto", "regression", "classification"), ml.options = ml_options())

ml_generalized_linear_regression() Same options for: **ml_gradient_boosted_trees**
ml_logistic_regression (x, response, features,
intercept = TRUE, family = gaussian(link = "identity"), iter.max = 100L, ml.options = ml_options())

ml_kmeans() (x, centers, iter.max = 100, features = dplyr::tbl_vars(x),
compute.cost = TRUE, tolerance = 1e-04, ml.options = ml_options())

ml_lda() (x, features = dplyr::tbl_vars(x), k = length(features), alpha = (50/k) + 1, beta = 0.1 + 1, ml.options = ml_options())

ml_linear_regression() (x, response, features, intercept = TRUE, alpha = 0, lambda = 0, iter.max = 100L, ml.options = ml_options())

ml_multilayer_perceptron() Same options for: **ml_logistic_regression**
(x, response, features, layers, iter.max = 100, seed = sample(.Machine\$integer.max, 1), ml.options = ml_options())

ml_naive_bayes() (x, response, features, lambda = 0, ml.options = ml_options())

ml_one_vs_rest() (x, classifier, response, features, ml.options = ml_options())

ml_pca() (x, features = dplyr::tbl_vars(x), ml.options = ml_options())

ml_random_forest() (x, response, features, max.bins = 32L, max.depth = 5L, num.trees = 20L, type = c("auto", "regression", "classification"), ml.options = ml_options())

ml_survival_regression() (x, response, features, intercept = TRUE, censor = "censor", iter.max = 100L, ml.options = ml_options())

ml_binary_classification_eval() (predicted_tbl_spark, label, score, metric = "areaUnderROC")

ml_classification_eval() (predicted_tbl_spark, label, predicted_lbl, metric = "f1")

ml_tree_feature_importance() (sc, model)

sparklyr

is an R
interface
for

Spark
Apache Spark