

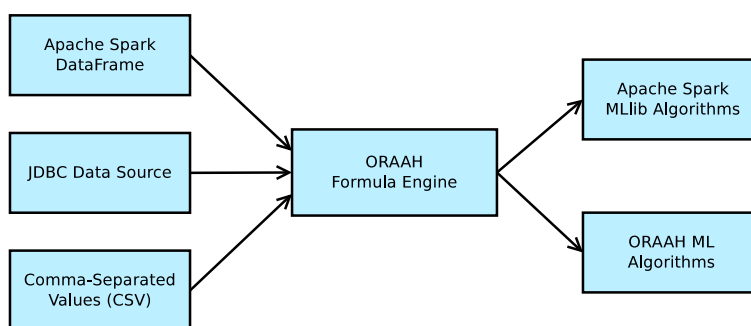
ORAAH 2.8.2 Formula and Data Preprocessing

May 15, 2020

1 Introduction

ORAAH 2.8.2 introduces new scalable data preprocessing algorithms to facilitate machine learning model building and deployment. ORAAH includes its own high performance Java implementation of the R formula to create parallel distributed model matrices, which serve as input for machine learning and statistical algorithms.

In this document we mainly use R to illustrate functionality. Please keep in mind, that both ORAAH Data Preprocessing and Formula engines are readily accessible from Java, Python, C++, R, Scala and other programming languages.



R formula is a remarkably powerful mechanism to define and create statistical and machine learning models. It is used to specify response and explanatory variables, nonlinear transformations, and interactions.

Here is a simple example

```
Y ~ X1 + X2
```

where Y is a dependent variable (target), and $X1$, $X2$ are predictors (features). The tilde operator \sim separates response (left-hand side) from explanatory terms (right-hand side of the formula).

Consider a table `tab` with three columns Y , $X1$, and $X2$. Y is a binary variable (only takes two values). We can fit (train) a logistic regression model as follows

```
fit <- orch.glm2(formula = Y ~ X1 + X2, data = tab)
```

Omitting the response (empty left-hand side) is used to define and create unsupervised models.

```
~ X1 + X2
```

With this specification, we can train a k-means model:

```
fit <- orch.ml.kmeans(formula = ~ X1 + X2, data = tab)
```

Operations allowed in a formula object fall into two main categories:

1. Set-theoretic. In our earlier example $X1 + X2$, the plus operator $+$ means to *include* variables as predictors to the algorithm for model building; there is no arithmetic summation here of any kind; the operation is precisely equivalent to how one adds elements (variables) to a set (model).
2. Arithmetic (e.g., add a number to each element of the column), where plus $+$, minus $-$, and multiply $*$ operators are used in their classical arithmetic sense. The divide operator $/$ can be used only in an arithmetic context. The places where these operators are understood in their arithmetic sense are *function arguments* (the most common case), response, and boolean expressions.

2 R Formula Syntax

$+$ plus operator when used in set-theoretic sense means to include the variable as a predictor to the algorithm for model building (into the set). To reiterate, there is no arithmetic summation here of any kind. For instance,

$$Y \sim A + B + C$$

here Y is our target (response); and we also include three predictors (data columns) A , B , and C into the model.

- minus operator removes the corresponding variable from the model.

() parenthesis are used to group variables.

$$(X + Y + Z)$$

creates a subset, which comprises three variables X , Y , and Z .

. dot-character stands for all variables in the training data, except the response variable. For instance

$$Y \sim . - (X20 + X50)$$

will include all variables (notice the dot-character), except $X20$ and $X50$. Alternatively, you can use $Y \sim . - X20 - X50$.

: colon operator generates and includes the interaction between the two variables.

$$A : B$$

* asterisk operator includes both the main effects, and the interactions between them; $A * B$ is equivalent to

$$A + B + A : B$$

It is perfectly fine to generate interactions between two groups of variables, for instance $(A + B) * (X + Y)$, is equivalent to $A * X + A * Y + B * X + B * Y$. Both colon and asterisk operators are of higher precedence than add (plus) and remove (minus).

$(A_1 + A_2 + \dots + A_k)^n$ include the variables and generate all interactions up to n-way. Example: $Y \sim (. - A)^3$. One more example

$$Y \sim (\log(A) + B : Z)^2$$

is equivalent to

$$Y \sim \log(A) + B : Z + \log(A) : B : Z$$

I() identity function. Its argument will be treated in the arithmetic sense. A handy function whenever you need to scale, or carry out any nonlinear transformation. Everything between the parentheses will be treated in the arithmetic sense. For instance

$$I(\log(A / 10) * B + C)$$

will create a new column, whose row elements will be

$$\log(A[\text{row}] / 10) * B[\text{row}] + C[\text{row}]$$

Here all arithmetic operators are understood in their traditional arithmetic sense. * means simple element-wise multiplication.

sin(), cos(), log(), exp(), ... math functions (see next sections). Example:

$$\log(A) \sim \exp(B)$$

relational operators traditional operators, & is a synonym for &&; similarly, | is a synonym for ||.

```
A >= B
A <= B
A > B
A < B
A == B
A != B
A && B
A & B
A || B
A | B
```

+1 adds intercept.

-1 removes intercept.

-0 adds intercept, synonym for +1.

+0 removes intercept, synonym for -1.

`as.factor(X)` allows treating a numerical column X (often integer) as a factor. For instance, consider the classical Airline on-time performance dataset, which has among others the following columns

- Cancelled, whether the flight was cancelled, binary integer column.
- Year, integer column.
- DayOfWeek, integer column.
- DepDelay departure delay in minutes, integer column.

We can train a logistic regression model as follows:

```
Cancelled ~ as.factor(Year) + as.factor(DayOfWeek) +
  DepDelay
```

The above formula will create three predictors: Year as a factor variable, DayOfWeek as a factor variable, and DepDelay as numeric. In ORAAH you can also use a shortcut `F(x)`, which means exactly the same thing as `as.factor(x)`, it is just a little less to type.

3 Arithmetic Functions in ORAAH Formula

Recall, functions treat their arguments in an arithmetic sense, similar to the identity function $I(x)$.

The argument x can be a number, an arithmetic expression that evaluates to number, or a column expression that is evaluated in the arithmetic sense. For instance

$$\log(Y) \sim I(X * Z * 0.1) + \tan(Z)$$

in the $I(X * Z * 0.1)$ expression above, the formula creates a new column whose row elements are the result of multiplication $X[\text{row}] * Z[\text{row}] * 0.1$.

Table 1: Arithmetic functions in ORAAH formula

<code>abs(x)</code>	absolute value
<code>acos(x)</code>	arc cosine
<code>asin(x)</code>	arc sine
<code>atan(x)</code>	arc tangent
<code>cbirt(x)</code>	cube root
<code>ceil(x)</code>	ceiling function. The smallest integer value (returned as double) that is greater than or equal to the argument. Returns the argument, if is already a whole number (integer).
<code>cos(x)</code>	trigonometric cosine
<code>cosh(x)</code>	hyperbolic cosine
<code>exp(x)</code>	exponent e^x
<code>expm1(x)</code>	$e^x - 1$
<code>floor(x)</code>	ceiling function, the largest integer value (returned as double) that is less than or equal to the argument. Returns the argument, if is already a whole number (integer).
<code>log(x)</code>	natural logarithm, $\ln(x)$
<code>log10(x)</code>	the base 10 logarithm
<code>log1p(x)</code>	$\ln(1 + x)$
<code>log2(x)</code>	base 2 logarithm
<code>rint(x)</code>	integer, closest to the argument
<code>round(x)</code>	integer, closest to the argument

<code>signum(x)</code>	signum function	$\begin{cases} +1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x > 0 \end{cases}$
<code>sin(x)</code>	trigonometric sine	
<code>sinh(x)</code>	hyperbolic sine	
<code>sqrt(x)</code>	square root	
<code>tan(x)</code>	trigonometric tangent	
<code>tanh(x)</code>	hyperbolic tangent	
<code>toDegrees(x)</code>	converts an angle measured in radians to an approximately equivalent angle measured in degrees	
<code>toRadians(x)</code>	converts an angle measured in degrees to an approximately equivalent angle measured in radians.	

4 Statistical Functions in ORAAH Formula

Statistical and special functions (described in the next section) can be very helpful to create new features, add noise to a dataset, etc.

Arguments `x`, `p`, `q` below must be column expressions, evaluated in an arithmetic sense. The remaining parameters (`alpha`, `beta`) must be numbers (or expressions which evaluate to numbers).

Table 2: Statistical functions in ORAAH formula

	Beta distribution	
	<code>alpha > 0, beta > 0</code>	
density	<code>dbeta(x, alpha, beta)</code>	
cumulative density	<code>pbeta(q, alpha, beta)</code>	
quantile	<code>qbeta(p, alpha, beta)</code>	
random deviates	<code>rbeta(alpha, beta)</code>	
	Binomial distribution	
	<code>nTrials > 0 integer, 0 <= probability <= 1</code>	
density	<code>dbinom(x, nTrials, probability)</code>	
cumulative density	<code>pbinom(q, nTrials, probability)</code>	
quantile	<code>qbinom(p, nTrials, probability)</code>	

random deviates	<code>rbinom(nTrials, probability)</code>
	Cauchy distribution
	<code>scale > 0</code>
density	<code>dcauchy(x, median, scale)</code>
cumulative density	<code>pcauchy(q, median, scale)</code>
quantile	<code>qcauchy(p, median, scale)</code>
random deviates	<code>rcauchy(median, scale)</code>
	Chi-squared distribution
	<code>degreesOfFreedom > 0</code>
density	<code>dchisq(x, degreesOfFreedom)</code>
cumulative density	<code>pchisq(q, degreesOfFreedom)</code>
quantile	<code>qchisq(p, degreesOfFreedom)</code>
random deviates	<code>rchisq(degreesOfFreedom)</code>
	Exponential distribution
	<code>rate > 0</code>
density	<code>dexp(x, rate)</code>
cumulative density	<code>pexp(q, rate)</code>
quantile	<code>qexp(p, rate)</code>
random deviates	<code>rexp(rate)</code>
	F-distribution
	<code>numeratorDF > 0, denominatorDF > 0</code>
density	<code>df(x, numeratorDF, denominatorDF)</code>
cumulative density	<code>pf(q, numeratorDF, denominatorDF)</code>
quantile	<code>qf(p, numeratorDF, denominatorDF)</code>
random deviates	<code>rf(numeratorDF, denominatorDF)</code>
	Gamma distribution
	<code>shape > 0, scale > 0</code>
density	<code>dgamma(x, shape, scale)</code>
cumulative density	<code>pgamma(q, shape, scale)</code>
quantile	<code>qgamma(p, shape, scale)</code>
random deviates	<code>rgamma(shape, scale)</code>

Geometric distribution

	$0 < \text{probability} \leq 1$
density	<code>dgeom(x, probability)</code>
cumulative density	<code>pgeom(q, probability)</code>
quantile	<code>qgeom(p, probability)</code>
random deviates	<code>rgeom(probability)</code>
Hypergeometric distribution	
	$\text{populationSize} > 0$
	$0 \leq \text{nSuccesses} \leq \text{populationSize}$
	$0 < \text{sampleSize} \leq \text{populationSize}$
density	<code>dhyper(x, populationSize, nSuccesses, sampleSize)</code>
cumulative density	<code>phyper(q, populationSize, nSuccesses, sampleSize)</code>
quantile	<code>qhyper(p, populationSize, nSuccesses, sampleSize)</code>
random deviates	<code>rhyper(populationSize, nSuccesses, sampleSize)</code>
Log-normal distribution	
	$\text{shape} > 0$
density	<code>dlnorm(x, scale, shape)</code>
cumulative density	<code>plnorm(q, scale, shape)</code>
quantile	<code>qlnorm(p, scale, shape)</code>
random deviates	<code>rlnorm(scale, shape)</code>
Normal distribution	
	$\text{sd} > 0$
density	<code>dnorm(x, mean, sd)</code>
cumulative density	<code>pnorm(q, mean, sd)</code>
quantile	<code>qnorm(p, mean, sd)</code>
random deviates	<code>rnorm(mean, sd)</code>
Poisson distribution	
	$\text{mean} > 0$
density	<code>dpois(x, mean)</code>
cumulative density	<code>ppois(q, mean)</code>
quantile	<code>qpois(p, mean)</code>
random deviates	<code>rpois(mean)</code>
Student t-distribution	
	$\text{degreesOfFreedom} > 0$

density	<code>dt(x, degreesOfFreedom)</code>
cumulative density	<code>pt(q, degreesOfFreedom)</code>
quantile	<code>qt(p, degreesOfFreedom)</code>
random deviates	<code>rt(degreesOfFreedom)</code>
Triangular distribution	
<code>lower <= mode <= upper</code>	
density	<code>dtriangular(x, lower, mode, upper)</code>
cumulative density	<code>ptriangular(q, lower, mode, upper)</code>
quantile	<code>qtriangular(p, lower, mode, upper)</code>
random deviates	<code>rtriangular(lower, mode, upper)</code>
Uniform distribution	
<code>lower < upper</code>	
density	<code>dunif(x, lower, upper)</code>
cumulative density	<code>punif(q, lower, upper)</code>
quantile	<code>qunif(p, lower, upper)</code>
random deviates	<code>runif(lower, upper)</code>
Weibull distribution	
<code>alpha > 0, beta > 0</code>	
density	<code>dweibull(x, alpha, beta)</code>
cumulative density	<code>pweibull(q, alpha, beta)</code>
quantile	<code>qweibull(p, alpha, beta)</code>
random deviates	<code>rweibull(alpha, beta)</code>
Pareto distribution	
<code>scale > 0, shape > 0</code>	
density	<code>dpareto(x, scale, shape)</code>
cumulative density	<code>ppareto(q, scale, shape)</code>
quantile	<code>qpareto(p, scale, shape)</code>
random deviates	<code>rpareto(scale, shape)</code>

5 Special Functions

Arguments of the special functions are treated in a numerical sense, and they can be numbers, expressions that evaluate to numbers, or column expressions.

Table 3: Special functions in ORAAH formula

<code>gamma(x)</code>	gamma function $\Gamma(x)$
<code>lgamma(x)</code>	natural logarithm of the gamma function $\ln(\Gamma(x))$
<code>digamma(x)</code>	digamma function, $\frac{d}{dx} \ln(\Gamma(x))$
<code>trigamma(x)</code>	trigamma function, $\frac{d^2}{dx^2} \ln(\Gamma(x))$
<code>lanczos(x)</code>	Lanczos approximation of the gamma function
<code>factorial(x)</code>	factorial $n!$
<code>lfactorial(x)</code>	natural logarithm of the factorial function $\ln(n!)$
<code>lbeta(a, b)</code>	natural logarithm of the Beta function $\ln\left(\frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}\right)$
<code>lchoose(n, k)</code>	natural logarithm of the binomial coefficient $\ln\left(\frac{n!}{k!(n-k)!}\right)$

6 Scaling and Aggregate Functions

Normalization can often be of paramount importance for successful creation and deployment of machine learning models in practice. In this section we describe new functions and techniques introduced in ORAAH 2.8.2

Notation:

x	data column
m	number of observations (e.g. number of rows in a table)
$\text{mean}(x)$	$= \frac{1}{m} \sum_{j=1}^m x_j$
$\text{sd}(x)$	$= \sqrt{\frac{\sum_{j=1}^m (x_j - \text{mean}(x))^2}{m - 1}}$
$\text{range}(x)$	$= \max(x) - \min(x)$
$\text{midrange}(x)$	$= \frac{\min(x) + \max(x)}{2}$

Table 4 lists normalization techniques for the function `orch.df.scale(data, method)`.

Table 4: Scalable (Parallel-Distributed) Normalization Techniques

method	what it is
"standardization"	$\frac{x - \text{mean}(x)}{\text{sd}(x)}$
"unitization"	$\frac{x - \text{mean}(x)}{\text{range}(x)}$
"unitization_zero_minimum"	$\frac{x - \min(x)}{\text{range}(x)}$
"normalization"	$\frac{x - \text{midrange}(x)}{\text{range}(x)/2}$
"normalization_2"	$\frac{x - \text{mean}(x)}{\max(x - \text{mean}(x))}$
"normalization_3"	$\frac{x - \text{mean}(x)}{\sqrt{\sum (x_j - \text{mean}(x))^2}}$
"quotient_sd"	$\frac{x}{\text{sd}(x)}$
"quotient_range"	$\frac{x}{\text{range}(x)}$
"quotient_max"	$\frac{x}{\max(x)}$
"quotient_mean"	$\frac{x}{\text{mean}(x)}$
"quotient_sum"	$\frac{x}{\sum x_j}$
"quotient_sqrt_ssq"	$\frac{x}{\sqrt{\sum x_j^2}}$

Consider an example where we create a Spark data frame from Comma-Separated Values (CSV) files; then create a summary (which would show us some basic statistics on the data frame columns); and then scale the columns.

```
# Create a CSV file; use a few columns ("GNP", "Unemployed", and "Employed")
# from the Longley's Economic Regression Data.
require(stats)
head(longley)
data <- subset(longley, select=c("GNP", "Unemployed", "Employed"))
write.csv(data, file="longley.csv", row.names=FALSE)
```

Now, our CSV dataset is ready. First, create a Spark data frame using the ORCH function:

```
orch.df.fromCSV(path, minPartitions= -1L,
  headerPresent=TRUE, fieldSeparator="," ,
  quote="\\"", na="NA", verbose=TRUE)
```

`orch.df.fromCSV()` can work in a fully automated “discovery” mode, where both CSV column names, and columns types will be automatically discovered. The CSV data can be read from HDFS, a local file system (in which case the file must exist on every compute node), or any other Hadoop-compliant file system.

```
dataFrame <- orch.df.fromCSV("longley.csv")
```

produces the following data frame (recall, we saved just a few columns from the longley dataset):

```
+-----+-----+-----+
|   GNP|Unemployed|Employed|
+-----+-----+-----+
|234.289|    235.6|  60.323|
|259.426|    232.5|  61.122|
|258.054|    368.2|  60.171|
|284.599|    335.1|  61.187|
|328.975|    209.9|  63.221|
|346.999|    193.2|  63.639|
|365.385|    187.0|  64.989|
|363.112|    357.8|  63.761|
|397.469|    290.4|  66.019|
| 419.18|    282.2|  67.857|
|442.769|    293.6|  68.169|
```

```

|444.546|    468.1| 66.513|
|482.704|    381.3| 68.655|
|502.601|    393.1| 69.564|
|518.173|    480.6| 69.331|
|554.894|    400.7| 70.551|
+-----+-----+-----+

```

Its summary, calculated by `orch.df.summary(dataFrame)`, is as follows (some digits were omitted for compactness):

```

+-----+-----+-----+-----+-----+-----+-----+-----+
| col_name| col_type|num_missing|  min|  max|      mean|      sd|num_factor_levels|
+-----+-----+-----+-----+-----+-----+-----+-----+
|      GNP|DoubleType|          0|234.289|554.894|    387.6...| 99.394...|          null|
|Unemployed|DoubleType|          0| 187.0| 480.6|319.33124999...| 93.446...|          null|
|  Employed|DoubleType|          0| 60.171| 70.551| 65.31700000...|3.51196...|          null|
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Now, let us scale the `dataFrame` via technique "normalization_2" (normalization in range [-1,1]):

```
scaledFrame <- orch.df.scale(dataFrame, "normalization_2")
```

This produces:

```

+-----+-----+-----+-----+-----+
|          GNP|          Unemployed|          Employed|
+-----+-----+-----+-----+-----+
| -0.9175449109183148| -0.5192031934271204| -0.9541459686664142|
| -0.7672000116629891| -0.5384257644459941| -0.8014902560183439|
| -0.7754059710765353|  0.3030267798318026| -0.9831868551776869|
| -0.6166397956883576| 0.09777932798511839| -0.7890714558654975|
| -0.35122605302398496| -0.6785645080029448| -0.40045854031333833|
| -0.24342414889151134| -0.7821183583304262| -0.3205961024073378|
| -0.13345711552601774| -0.8205635003681736| -0.06266717615590434|
| -0.14705197394219113|  0.2385381544781617| -0.29728696981276387|
|  0.05843792953536062| -0.17939774444831974|  0.13412304165074493|
|  0.1882918543367441| -0.23024454520792126|  0.4852884982804728|
|  0.32937813466191723| -0.159555090493353|  0.5448987390141369|
|  0.34000640716765423|  0.9224896329884122|  0.22850592281238044|
|  0.568230167591918|  0.3842576444599467|  0.6377531524646537|
|  0.6872345221482776|  0.45742743091888566|  0.8114252961406171|
|  0.7803709652880291|                    1.0|  0.7669086740542606|
|                    1.0|  0.5045537340619308|                    1.0|
+-----+-----+-----+-----+-----+

```

Its summary is as follows (some digits were omitted for compactness):

```

+-----+-----+-----+-----+-----+-----+-----+-----+
|col_name |col_type |num_missing|min      |max|mean |sd      |num_factor_levels|
+-----+-----+-----+-----+-----+-----+-----+-----+
|GNP      |DoubleType|0          |-0.917...|1.0|-5e-17|0.594...|null              |
|Unemployed|DoubleType|0          |-0.820...|1.0| 3e-16|0.579...|null              |
|Employed |DoubleType|0          |-0.983...|1.0|-1e-15|0.670...|null              |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Note: `orch.df.scale()` preserves the original column names, but it changes integral column types (byte, short, integer, long), and 32-bit floating point into 64-bit double. In other words, the result of scaling is a 64-bit double column. Boolean, date, calendar, and timestamp columns are treated as strings (factor), and remain unchanged.

In addition to the `orch.df.scale()` users can use the following aggregate functions, which return a single double floating point value:

Table 5: Formula aggregate functions

function name	what it is
(colExpr below stands for column expression)	
<code>avg(colExpr)</code>	mean (average) value
<code>mean(colExpr)</code>	mean value (synonym for avg)
<code>max(colExpr)</code>	maximum value
<code>min(colExpr)</code>	minimum value
<code>sd(colExpr)</code>	standard deviation
<code>stddev(colExpr)</code>	standard deviation (synonym for sd)
<code>sum(colExpr)</code>	sum
<code>variance(colExpr)</code>	variance
<code>var(colExpr)</code>	variance (synonym for variance)
<code>kurtosis(colExpr)</code>	kurtosis
<code>skewness(colExpr)</code>	skewness

Each aggregate function returns a single double floating point value; and therefore, can only be used inside an arithmetic context. In fact, this is exactly how formula preprocessor works: it will compute each aggregate function value, and then literally replace aggregates with their values. For instance

```
Kyphosis ~ I( Start / max(Start) )
```

In the above example, `I()` creates an arithmetic context, where we use `max(Start)` as a denominator to scale the `Start` variable. Formula pre-

processor will compute the $\max(\text{Start})$ (18.0), and the final formula will become

$$\text{Kyphosis} \sim I(\text{Start} / 18.0)$$

This value 18.0 will remain fixed (will be used for prediction / scoring).