



## BigML Dataset Cheat Sheet



### Sampling

#### Missing

Operation	Description	Field Type	API Name	Operation	Description	Field Type	API Name
If value is missing	Includes instances containing missing values for the selected field.	All	json_filter or lisp_filter	Is like (case sensitive)	Matches words containing at least part of the letters specified taking into account lower and upper cases. e.g., "great" will also match text containing the word "great" or "greatness", but not "Great" or "Greatness".	Text	json_filter or lisp_filter
If value is not missing	Excludes instances containing missing values for the selected field.	All	json_filter or lisp_filter	Is like (case insensitive)	Matches words containing at least part of the letters specified not taking into account lower and upper cases. e.g., "great" will also match text containing the words "great", "greatness", "Great" or "Greatness".	Text	json_filter or lisp_filter
API Name				Contains (case sensitive)	Matches texts containing the exact words specified taking into account lower and upper cases. e.g., "great" will match text containing the word "great", but not "Great".	Text	json_filter or lisp_filter
Default				Contains (case insensitive)	Matches texts containing the exact words specified not taking into account lower and upper cases. e.g., "great" will match text containing the word "great" or "Great".	Text, items	json_filter or lisp_filter
Option	Description						
Rate	Sets the proportion of the dataset you want to consider between 0% and 100%.	100%	sample_rate				
Range	Specifies a subset of instances from which to sample, e.g., from instance 5 to instance 1000. The Rate you set will be computed over the Range configured.	(1, max_rows_in_dataset)					
Sampling	Allows you to choose between a random sampling or a deterministic sampling. When using deterministic sampling the random number generator will always use the same seed, producing repeatable results.	Random seed					
Replacement	Allows a single instance to be selected multiple times. Sampling without replacement ensures that each instance cannot be selected more than once.	False	replacement				

### Comparison

Operation	Description	Field Type	API Name	Operation	Description	Field Type	API Name
is between	Includes instances containing values within the specified range.	Numeric	json_filter or lisp_filter	Is not like (case sensitive)	Checks whether given words do not match a regular expression taking into account lower and upper cases. It is a combination of "not" and "matches"; e.g., it will filter instances that do not contain "great" or "greatness", but will not consider "Great" or "Greatness".	Text	json_filter or lisp_filter
is less than	Includes instances containing values below the specified level.	Numeric	json_filter or lisp_filter	Is not like (case insensitive)	Same behavior as above but not taking into account lower and upper cases, e.g., it will filter instances that do not contain "great", "greatness", "Great" or "Greatness".	Text	json_filter or lisp_filter
is less than or equal to	Includes instances containing values equal or below the specified level.	Numeric	json_filter or lisp_filter	Not contains (case sensitive)	Checks whether a given text does not match a regular expression taking into account lower and upper cases. It is a combination of "not" and "matches"; e.g., it will filter instances that do not contain "great" or "greatness", but will not consider "Great" or "Greatness".	Text	json_filter or lisp_filter
is greater than or equal to	Includes instances containing values equal or above the specified level.	Numeric	json_filter or lisp_filter	Not contains (case insensitive)	Same behavior as above but not taking into account lower and upper cases, e.g., it will filter instances that do not contain "great", "greatness", "Great" or "Greatness".	Text, items	json_filter or lisp_filter
is between percentiles	Includes instances within the specified percentiles. e.g., a percentile between 0 and 0.5 includes the first 30% of the instances.	Numeric	json_filter or lisp_filter				
API Name							
JSON flattening formula	Computes any operation using the <b>Flattine JSON syntax</b> .	All	json_filter or lisp_filter	Operation	Description	Field Type	API Name
Lisp flattening formula	Computes any operation using the <b>Flattine Lisp-like syntax</b> .	All	lisp_filter   field   fields				



### Pre-Defined Operations to Filter Datasets

#### Equality

Operation	Description	Field Type	API Name	Operation	Description	Field Type	API Name
Equals	Includes instances containing the specified value/values.	All	json_filter or lisp_filter	Is below the mean	Includes instances below the mean of the selected field.	Numeric	json_filter or lisp_filter
Does not equal	Excludes instances containing the specified value/values.	All	json_filter or lisp_filter	Is above the mean	Includes instances above the mean of the selected field.	Numeric	json_filter or lisp_filter

### Flattine Formula

# Pre-Defined Operations to Add New Fields to Your Dataset

## Normalization

## Sliding Windows

Normalization				Sliding Windows			
Operation	Description	Field Type	API Name	Operation	Description	Field Type	API Name
<b>Discretization</b>	Normalizes your field. Select the range for which you want to normalize your field. (This range should be within the field range.)	Numeric	new_fields (+ Flatline)	Sum of Instances	Sums consecutive instances by defining a window start and end (negative values are previous instances and positive values next instances).	Numeric	new_fields (+ Flatline)
<b>Operation</b>	<b>API Name</b>	<b>Field Type</b>	<b>API Name</b>	<b>Operation</b>	<b>Description</b>	<b>Field Type</b>	<b>API Name</b>
<b>Discretize by percentiles</b>	Splits your field values into equal population segments (Categories), e.g., setting 3 groups for a field ranging from 0 to 6 will yield: category 1 = [0, 2], category 2 = [2, 4], category 3 = [4, 6].	Numeric	new_fields (+ Flatline)	Z-score	Indicates the distance of the values from the mean.	Numeric	new_fields (+ Flatline)
<b>Is within percentiles?</b>	Computes a boolean field with True or False values for each instance when you specify a percentile range between 0 and 1, depending whether they belong to the specified range or not.	Numeric	new_fields (+ Flatline)				
<b>Maths</b>				<b>Field Type</b>	<b>API Name</b>	<b>Field Type</b>	<b>API Name</b>
<b>Fixed value</b>	Replaces all your field missing values by the specified value. You can set a number or a string.	Numeric, categorical	new_fields (+ Flatline)	<b>Exponentiation</b>	Computes $e^x$ elevated to the field value: $e^x$	Numeric	new_fields (+ Flatline)
<b>Maximum</b>	Replaces missing values by the max. value of the selected field.	Numeric	new_fields (+ Flatline)	<b>log2</b>	Scales fields logarithmically with a logarithm base of 2. This is useful for fields with a wide range of data (since it reduces the range into a more manageable scale) and to find exponential patterns in your data.	Numeric	new_fields (+ Flatline)
<b>Mean</b>	Replaces missing values by the mean of the selected field.	Numeric	new_fields (+ Flatline)	<b>log</b>	Scales fields logarithmically, with a logarithm base of e. This is useful for fields with a wide range of data (since it reduces the range into a more manageable scale) and to find exponential patterns in your data.	Numeric	new_fields (+ Flatline)
<b>Minimum</b>	Replaces missing values by the min. value of the selected field.	Numeric	new_fields (+ Flatline)	<b>ln</b>	Scales fields logarithmically, with a logarithm base of e. This is useful for fields with a wide range of data (since it reduces the range into a more manageable scale) and to find exponential patterns in your data.	Numeric	new_fields (+ Flatline)
<b>Population</b>	Replaces missing values by the number of the total instances that have valid values for the selected field, e.g., for a field containing 54 instances with valid values, the missing values will be replaced by 54.	Numeric	new_fields (+ Flatline)	<b>Square</b>	Squares field values: $x^2$	Numeric	new_fields (+ Flatline)
<b>Random integer</b>	Replaces missing values by a random value. You can set the max. value you want for your random value generator.	Numeric	new_fields (+ Flatline)	<b>Square root</b>	Computes the square root of the value: $\sqrt{x}$	Numeric	new_fields (+ Flatline)
<b>Random value</b>	Replaces missing values by a random value within your field range.	Numeric, categorical	new_fields (+ Flatline)				
<b>Random weighted value</b>	Replaces missing values by a random value within the field range, using the same probability distribution as the values in the dataset (which is described by the field's histogram).	Numeric, categorical	new_fields (+ Flatline)				
<b>Replacing Missing Values</b>				<b>Operation</b>	<b>Description</b>	<b>Field Type</b>	<b>API Name</b>
<b>Types</b>				<b>Categorical</b>	Coerces numeric field values into categorical values, e.g., the number 10 will become a string "10".	Categorical	new_fields (+ Flatline)
<b>Integer</b>				<b>Boolean</b>	Coerces categorical values to integer values. Boolean values are assigned 0 (false) and 1 (true).	Text	new_fields (+ Flatline)
<b>Real</b>					Coerces categorical values to floating point values, e.g., the string "7.5 pounds" will become 7.5. Boolean values are assigned 0 and 1.	Text	new_fields (+ Flatline)

## Random

Operation	Description	Field Type	API Name
<b>Random integer</b>	Computes a random integer for each instance.	Numeric	new_fields (+ Flatline)
<b>Random value within field range</b>	Computes a random value by taking your field range as the reference for min. and max. values.	Numeric, categorical	new_fields (+ Flatline)
<b>Random weighted value</b>	Computes a random value within the field's range, using the same probability distribution as the values in the dataset (which is described by the field's histogram).	Numeric, categorical	new_fields (+ Flatline)

## Statistics

Operation	Description	Field Type	API Name
<b>Mean</b>	Computes the field mean for all instances.	Numeric	new_fields (+ Flatline)
<b>Population</b>	Computes the count of total instances for that field.	Numeric	new_fields (+ Flatline)
<b>Population fraction</b>	Computes the number of instances with values below the specified value.	Numeric	new_fields (+ Flatline)

## Flatline Formula

Operation	Description	Field Type	API Name
<b>JSON flatline formula</b>	Computes any operation using the <b>Flatline</b> JSON syntax.	All	new_fields["field"]
<b>Lisp flatline formula</b>	Computes any operation using the <b>Flatline</b> Lisp-like syntax.	All	new_fields["field"]

## Remove Duplicates

## Join Datasets

## Merge Datasets

Option	Description	Field Type	API Name
Remove duplicates	Removes the duplicated instances in your dataset taking into account all the field values.	All	SQL query

Option	Description	Field Type	API Name
Merge datasets	Select up to 32 datasets with the same fields to merge their instances into one dataset.	N/A	SQL query
Rate	Sets the proportion of each merging dataset you want to consider between 0% and 100%.	100%	sample_rate
Type of join	Allows you to select a left join, right join, full join, and inner join.	N/A	SQL query
Selected dataset	Allows you to select the dataset you want to join with the current dataset.	N/A	SQL query
Join fields	Allows you to select one or more fields from the current dataset to match the instances with the selected dataset. These fields should have the same values in both datasets so the instances can be matched.	All	SQL query
Select fields	Allows you to choose all the fields from the selected dataset or select a subset of them.	All	SQL query
Filter datasets	Allows you to filter the current and/or selected dataset while making the join.	All	SQL query

## Aggregate Instances

Option	Description	Field Type	API Name
Aggregating field	Sets the field of the dataset that you want to use to group your instances.	All	SQL query
Operation	Specifies the aggregation operation to be applied to the rest of the dataset fields.	All	SQL query

## Order Instances

Option	Description	Default	API Name
Order Instances	Allows you to sort the rows of a dataset by one or more selected fields in ascending or descending order	N/A	SQL query