

Package ‘mlr3oml’

October 12, 2022

Title Connector Between 'mlr3' and 'OpenML'

Version 0.6.0

Description Provides an interface to 'OpenML.org' to list and download machine learning data, tasks and experiments. The 'OpenML' objects can be automatically converted to 'mlr3' objects.
For a more sophisticated interface which also allows uploading to 'OpenML', see the 'OpenML' package.

License LGPL-3

URL <https://mlr3oml.mlr-org.com>, <https://github.com/mlr-org/mlr3oml>

BugReports <https://github.com/mlr-org/mlr3oml/issues>

Depends R (>= 3.1.0)

Imports backports (>= 1.1.6), checkmate, curl, data.table, jsonlite, lgr, methods, mlr3 (>= 0.14.0), mlr3misc (>= 0.7.0), paradox, R6, stringi, uuid, withr

Suggests DBI, duckdb, mlr3db (>= 0.5.0), qs, RWeka, testthat (>= 3.0.0)

Config/testthat/edition 3

Encoding UTF-8

NeedsCompilation yes

RoxygenNote 7.2.1

Author Michel Lang [aut] (<<https://orcid.org/0000-0001-9754-0393>>),
Sebastian Fischer [cre, aut] (<<https://orcid.org/0000-0002-9609-3197>>)

Maintainer Sebastian Fischer <sebf.fischer@gmail.com>

Repository CRAN

Date/Publication 2022-10-12 12:10:12 UTC

R topics documented:

mlr3oml-package	2
as_learner.OMLFlow	3

benchmark_grid_oml	4
list_oml	4
oml_collection_r6	8
oml_data_r6	11
oml_flow_r6	14
oml_object_r6	16
oml_run_r6	18
oml_sugar	20
oml_task_r6	21
read_arff	24
write_arff	25

Index	26
--------------	-----------

mlr3oml-package	<i>mlr3oml: Connector Between 'mlr3' and 'OpenML'</i>
-----------------	---

Description

Provides an interface to 'OpenML.org' to list and download machine learning data, tasks and experiments. The 'OpenML' objects can be automatically converted to 'mlr3' objects. For a more sophisticated interface which also allows uploading to 'OpenML', see the 'OpenML' package.

mlr3 Integration

This package adds the `mlr3::Task "oml"` and the `mlr3::Resampling "oml"` to `mlr3::mlr_tasks` and `mlr3::mlr_resamplings`, respectively. For the former you may pass either a `data_id` or a `task_id`, the latter requires a `task_id`. Furthermore it allows to convert the OpenML objects to mlr3 objects using the usual S3 generics such as `mlr3::as_task`, `mlr3::as_learner`, `mlr3::as_resampling`, `mlr3::as_resample_result`, `mlr3::as_benchmark_result` or `mlr3::as_data_backend`. This allows for a frictionless integration of OpenML and mlr3.

Options

- `mlr3oml.cache`: Enables or disables caching globally. If set to `FALSE`, caching is disabled. If set to `TRUE`, cache directory as reported by `R_user_dir()` is used. Alternatively, you can specify a path on the local file system here. Default is `FALSE`.
- `mlr3oml.api_key`: API key to use. All operations supported by this package work without an API key, but you might get rate limited without an API key. If not set, defaults to the value of the environment variable `OPENMLAPIKEY`.
- `mlr3oml.arff_parser`: ARFF parser to use, defaults to the internal one relies on `data.table::fread()`. Can also be set to "RWeka" for the parser in **RWeka**.
- `mlr3oml.parquet`: Enables or disables parquet as the default file format. If set to `TRUE`, the parquet version of datasets will be used by default. If set to `FALSE`, the arff version of datasets will be used by default. Note that the OpenML sever is still transitioning from arff to parquet and some features will work better with arff. Default is `FALSE`.

Relevant for developers

- `mlr3oml.test_server`: The default value for whether to use the OpenML [test server](#). Default is FALSE.
- `mlr3oml.test_api_key`: API key to use for the test server. If not set, defaults to the value of the environment variable TESTOPENMLAPIKEY.

Logging

The `lgr` package is used for logging. To change the threshold, use `lgr::get_logger("mlr3oml")$set_threshold()`.

Author(s)

Maintainer: Sebastian Fischer <sebf.fischer@gmail.com> ([ORCID](#))

Authors:

- Michel Lang <michellang@gmail.com> ([ORCID](#))

See Also

Useful links:

- <https://mlr3oml.ml-org.com>
- <https://github.com/mlr-org/mlr3oml>
- Report bugs at <https://github.com/mlr-org/mlr3oml/issues>

as_learner.OMLFlow *Convert an OpenML Flow to a mlr3 Learner*

Description

By default this function creates a Pseudo-Learner (that cannot be used for training or prediction) for the given task type. This enables the conversion of OpenML Runs to `mlr3::ResampleResults`. This is well defined because each subcomponent (i.e. id) can only appear once in a Flow according to the OpenML docs.

Usage

```
## S3 method for class 'OMLFlow'
as_learner(x, task_type = NULL, ...)
```

Arguments

<code>x</code>	(OMLFlow) The OMLFlow that is converted to a <code>mlr3::Learner</code> .
<code>task_type</code>	(character(1)) The task type to construct a pseudo-learner. For more information see OMLFlow .
<code>...</code>	Additional arguments.

benchmark_grid_oml *Helper function to create a benchmark design*

Description

OMLTasks contain tasks as well as resamplings. In order to create a benchmark design from a list of tasks and corresponding instantiated resamplings, this function can be used.

Usage

```
benchmark_grid_oml(tasks, learners, resamplings)
```

Arguments

tasks (list() or Task) A list of **mlr3::Tasks**.

learners (list() or Learner) A list of **mlr3::Learners**.

resamplings (list() or Resampling) A list of **mlr3::Resamplings** that are instantiated on the given tasks.

Value

([data.table\(\)](#))

Examples

```
library("mlr3")
collection = OMLCollection$new(258)
otasks = collection$tasks[1:2, ][["task"]]
tasks = as_tasks(otasks)
resamplings = as_resamplings(otasks)
learners = lrns(c("classif.rpart", "classif.featureless"))
design = benchmark_grid_oml(tasks, learners, resamplings)
print(design)
bmr = benchmark(design)
```

list_oml *List Data from OpenML*

Description

This function allows to query data sets, tasks, flows, setups, runs, and evaluation measures from <https://www.openml.org/d> using some simple filter criteria.

Usage

```
list_oml_data(  
  data_id = NULL,  
  data_name = NULL,  
  number_instances = NULL,  
  number_features = NULL,  
  number_classes = NULL,  
  number_missing_values = NULL,  
  tag = NULL,  
  limit = limit_default(),  
  test_server = test_server_default(),  
  ...  
)  
  
list_oml_evaluations(  
  run_id = NULL,  
  task_id = NULL,  
  measures = NULL,  
  tag = NULL,  
  limit = limit_default(),  
  test_server = test_server_default(),  
  ...  
)  
  
list_oml_flows(  
  uploader = NULL,  
  tag = NULL,  
  limit = limit_default(),  
  test_server = test_server_default(),  
  ...  
)  
  
list_oml_measures(test_server = test_server_default())  
  
list_oml_runs(  
  run_id = NULL,  
  task_id = NULL,  
  tag = NULL,  
  flow_id = NULL,  
  limit = limit_default(),  
  test_server = test_server_default(),  
  ...  
)  
  
list_oml_setups(  
  flow_id = NULL,  
  setup_id = NULL,  
  tag = NULL,
```

```

    limit = limit_default(),
    test_server = test_server_default(),
    ...
)

list_oml_tasks(
  task_id = NULL,
  data_id = NULL,
  number_instances = NULL,
  number_features = NULL,
  number_classes = NULL,
  number_missing_values = NULL,
  tag = NULL,
  limit = limit_default(),
  test_server = test_server_default(),
  ...
)

```

Arguments

<code>data_id</code>	(integer()) Vector of data ids to restrict to.
<code>data_name</code>	(character(1)) Filter for name of data set.
<code>number_instances</code>	(integer()) Filter for number of instances.
<code>number_features</code>	(integer()) Filter for number of features.
<code>number_classes</code>	(integer()) Filter for number of labels of the target (only classification tasks).
<code>number_missing_values</code>	(integer()) Filter for number of missing values.
<code>tag</code>	(character()) Filter for tags. You can provide multiple tags as character vector.
<code>limit</code>	(integer()) Limit the results to <code>limit</code> records. Default is the value of option <code>"mlr3oml.limit"</code> , defaulting to 5000.
<code>test_server</code>	(character(1)) Whether to use the OpenML test server (https://test.openml.org/) or public server (https://www.openml.org/). Defaults to value of option <code>"mlr3oml.test_server"</code> , or FALSE if not set.
<code>...</code>	(any) Additional (unsupported) filters, as named arguments.

run_id	(integer()) Vector of run ids to restrict to.
task_id	(integer()) Vector of task ids to restrict to.
measures	(character()) Vector of evaluation measures to restrict to.
uploader	(integer(1)) Filter for uploader.
flow_id	(integer(1)) Filter for flow id.
setup_id	(integer()) Vector of setup ids to restrict to.

Details

Filter values are usually provided as single atomic values (typically integer or character). Provide a numeric vector of length 2 ($c(l, u)$) to find matches in the range $[l, u]$.

Note that only a subset of filters is exposed here. For a more feature-complete package, see **OpenML**. Alternatively, you can pass additional filters via `...` using the names of the official API, c.f. https://www.openml.org/api_docs.

Value

(`data.table()`) of results, or a null `data.table` if no data set matches the filter criteria.

References

Casalicchio G, Bossek J, Lang M, Kirchhoff D, Kerschke P, Hofner B, Seibold H, Vanschoren J, Bischl B (2017). “OpenML: An R Package to Connect to the Machine Learning Platform OpenML.” *Computational Statistics*, 1–15. doi: [10.1007/s0018001707422](https://doi.org/10.1007/s0018001707422).

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). “OpenML.” *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
### query data sets
# search for titanic data set
data_sets = list_oml_data(data_name = "titanic")
print(data_sets)

# search for a reduced version
data_sets = list_oml_data(
  data_name = "titanic",
  number_instances = c(2200, 2300),
  number_features = 4
)
print(data_sets)
```

```

### search tasks for this data set
tasks = list_oml_tasks(data_id = data_sets$data_id)
print(tasks)

# query runs, group by number of runs per task_id
runs = list_oml_runs(task_id = tasks$task_id)
runs[, .N, by = task_id]

```

oml_collection_r6 *OpenML Collection*

Description

This is the class for collections (previously known as studies) served on <https://www.openml.org>. A collection can either be a **task collection** or **run collection**. This object can also be constructed using the sugar function `oml_collection()`.

Run Collection

A run collection contains runs, flows, datasets and tasks. The primary object are the runs (`main_entity_type` is "run"). The the flows, datasets and tasks are those used in the runs.

Task Collection A task collection (`main_entity_type = "task"`) contains tasks and datasets. The primary object are the tasks (`main_entity_type` is "task"). The datasets are those used in the tasks.

Note: All Benchmark Suites on OpenML are also collections.

Caching

The OpenML collection itself cannot be not cached, this is because it can be modified in-place on the server, e.g. by adding or removing tasks or runs. The construction argument `cache` therefore only controls whether caching is applied to the OpenML objects that are contained in the collection.

mlr3 Intergration

- Obtain a list of `mlr3::Tasks` using `mlr3::as_tasks`.
- Obtain a list of `mlr3::Resamplings` using `mlr3::as_resamplings`.
- Obtain a list of `mlr3::Learners` using `mlr3::as_learners` (if `main_entity_type` is "run").
- Obtain a `mlr3::BenchmarkResult` using `mlr3::as_benchmark_result` (if `main_entity_type` is "run").

Super class

`mlr3oml::OMLObject` -> `OMLCollection`

Active bindings

`desc (list())`
Collection description (meta information), downloaded and converted from the JSON API response.

`parquet (logical(1))`
Whether to use parquet.

`main_entity_type (character(n))`
The main entity type, either "run" or "task".

`flow_ids (integer(n))`
An vector containing the flow ids of the collection.

`data_ids (integer(n))`
An vector containing the data ids of the collection.

`run_ids (integer(n))`
An vector containing the run ids of the collection.

`task_ids (integer(n))`
An vector containing the task ids of the collection.

`runs (data.table())` A data.table summarizing the runs included in the collection. Returns NULL for Task Collections.

`flows (data.table())` A data.table summarizing the flows included in the collection. Returns NULL for Task Collections.

`data (data.table())` A data.table summarizing the datasets included in the collection.

`tasks (data.table())` A data.table summarizing the tasks included in the collection.

Methods

Public methods:

- [OMLCollection\\$new\(\)](#)
- [OMLCollection\\$print\(\)](#)
- [OMLCollection\\$clone\(\)](#)

Method `new()`: Creates a new instance of this [R6](#) class.

Usage:

```
OMLCollection$new(  
  id,  
  cache = cache_default(),  
  parquet = parquet_default(),  
  test_server = test_server_default()  
)
```

Arguments:

`id (integer(1))`
OpenML id for the object.

cache (logical(1) | character(1))

See field cache for an explanation of possible values. Defaults to value of option "mlr3oml.cache", or FALSE if not set. The collection itself is not cached, this is because it can be modified in-place on OpenML, e.g. by adding or removing tasks or runs. This parameter therefore only controls whether the contained elements are cached when loaded, e.g. when accessing the included tasks.

parquet (logical(1))

Whether to use parquet instead of arff. If parquet is not available, it will fall back to arff. Defaults to value of option "mlr3oml.parquet" or FALSE if not set.

test_server (character(1))

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>). Defaults to value of option "mlr3oml.test_server", or FALSE if not set.

Method print(): Prints the object.

Usage:

```
OMLCollection$print()
```

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
OMLCollection$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

References

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). "OpenML." *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
library("mlr3")
# OpenML Run collection:
run_collection = OMLCollection$new(id = 232)
# using sugar
run_collection = oml_collection(id = 232)
run_collection$main_entity_type
run_collection$tasks
run_collection$data
run_collection$flows
run_collection$runs

# mlr3 conversion:
tasks = as_tasks(run_collection)
resamplings = as_resamplings(run_collection)
learners = as_learners(run_collection, "classif")

bmr = as_benchmark_result(run_collection)
bmr$score(msr("classif.ce"))
```

```

# OpenML task collection
task_collection = OMLCollection$new(id = 258)
# using sugar
task_collection = oml_collection(id = 258)

task_collection$main_entity_type
task_collection$tasks
task_collection$data

# mlr3 conversion
tasks = as_tasks(task_collection)
resamplings = as_resamplings(task_collection)

```

oml_data_r6

Interface to OpenML Data Sets

Description

This is the class for data sets served on **OpenML**. This object can also be constructed using the `sugar` function `oml_data()`.

mlr3 Integration

- A `mlr3::Task` can be obtained by calling `mlr3::as_task()`.
- A `mlr3::DataBackend` can be obtained by calling `mlr3::as_data_backend()`. Depending on the selected file-type, the returned backend is a `mlr3::DataBackendDataTable` (arff) or `mlr3db::DataBackendDuckDB` (parquet).

Name conversion

Column names that don't comply with R's naming scheme are renamed (see `base::make.names()`). This means that the names can differ from those on OpenML.

File Format

The datasets stored on OpenML are either stored as (sparse) ARFF or parquet. When creating a new `OMLData` object, the constructor argument `parquet` allows to switch between arff and parquet. Note that not necessarily all data files are available as parquet. The option `mlr3oml.parquet` can be used to set a default. If `parquet` is TRUE but not available, "arff" will be used as a fallback.

ARFF Files

This package comes with an own reader for ARFF files, based on `data.table::fread()`. For sparse ARFF files and if the **RWeka** package is installed, the reader automatically falls back to the implementation in (`RWeka::read.arff()`).

Parquet Files

For the handling of parquet files, we rely on **duckdb** and CRANpkg{DBI}.

Super class

```
mlr3oml::OMLObject -> OMLData
```

Active bindings

qualities (data.table())

Data set qualities (performance values), downloaded from the JSON API response and converted to a `data.table::data.table()` with columns "name" and "value".

tags (character())

Returns all tags of the object.

parquet (logical(1))

Whether to use parquet.

data (data.table())

Returns the data (without the row identifier and ignore id columns).

features (data.table())

Information about data set features (including target), downloaded from the JSON API response and converted to a `data.table::data.table()` with columns:

- "index" (integer()): Column position.
- "name" (character()): Name of the feature.
- "data_type" (factor()): Type of the feature: "nominal" or "numeric".
- "nominal_value" (list()): Levels of the feature, or NULL for numeric features.
- "is_target" (logical()): TRUE for target column, FALSE otherwise.
- "is_ignore" (logical()): TRUE if this feature should be ignored. Ignored features are removed automatically from the data set.
- "is_row_identifier" (logical()): TRUE if the column encodes a row identifier. Row identifiers are removed automatically from the data set.
- "number_of_missing_values" (integer()): Number of missing values in the column.

target_names (character())

Name of the default target, as extracted from the OpenML data set description.

feature_names (character())

Name of the features, as extracted from the OpenML data set description.

nrow (integer())

Number of observations, as extracted from the OpenML data set qualities.

ncol (integer())

Number of features (including targets), as extracted from the table of data set features. This excludes row identifiers and ignored columns.

license (character())

Returns all license of the dataset.

parquet_path (character())

Downloads the parquet file (or loads from cache) and returns the path of the parquet file. Note that this also normalizes the names of the parquet file.

Methods

Public methods:

- [OMLData\\$new\(\)](#)
- [OMLData\\$print\(\)](#)
- [OMLData\\$quality\(\)](#)
- [OMLData\\$clone\(\)](#)

Method `new()`: Creates a new instance of this [R6](#) class.

Usage:

```
OMLData$new(
  id,
  cache = cache_default(),
  parquet = parquet_default(),
  test_server = test_server_default()
)
```

Arguments:

`id` (`integer(1)`)

OpenML id for the object.

`cache` (`logical(1)` | `character(1)`)

See field `cache` for an explanation of possible values. Defaults to value of option `"mlr3oml.cache"`, or `FALSE` if not set.

`parquet` (`logical(1)`)

Whether to use `parquet` instead of `arff`. If `parquet` is not available, it will fall back to `arff`.

Defaults to value of option `"mlr3oml.parquet"` or `FALSE` if not set.

`test_server` (`character(1)`)

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>).

Defaults to value of option `"mlr3oml.test_server"`, or `FALSE` if not set.

Method `print()`: Prints the object. For a more detailed printer, convert to a [mlr3::Task](#) via `as_task()`.

Usage:

```
OMLData$print()
```

Method `quality()`: Returns the value of a single OpenML data set quality.

Usage:

```
OMLData$quality(name)
```

Arguments:

`name` (`character(1)`)

Name of the quality to extract.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
OMLData$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

References

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). “OpenML.” *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
library("mlr3")
# OpenML Data object
odata = OMLData$new(id = 9)
# using sugar
odata = oml_data(id = 9)
print(odata)
print(odata$target_names)
print(odata$feature_names)
print(odata$tags)

# mlr3 conversion:
task = as_task(odata)
backend = as_data_backend(odata)
class(backend)

# get a task via tsk():
tsk("oml", data_id = 9)

# For parquet files
if (requireNamespace("duckdb")) {
  odata = OMLData$new(id = 9, parquet = TRUE)
  # using sugar
  odata = oml_data(id = 9)

  print(odata)
  print(odata$target_names)
  print(odata$feature_names)
  print(odata$tags)

  backend = as_data_backend(odata)
  class(backend)
  task = as_task(odata)
  task = tsk("oml", data_id = 9, parquet = TRUE)
  class(task$backend)
}
```

Description

This is the class for flows served on [OpenML](#). Flows represent machine learning algorithms. This object can also be constructed using the sugar function `oml_flow()`.

mlr3 Integration

- Obtain a `mlr3::Learner` using `mlr3::as_learner()`.

Super class

`mlr3oml::OMLObject` -> `OMLFlow`

Active bindings

`parameter` (`data.table`)
The parameters of the flow.

`dependencies` (`character()`)
The dependencies of the flow.

`tags` (`character()`)
Returns all tags of the object.

Methods**Public methods:**

- `OMLFlow$new()`
- `OMLFlow$print()`
- `OMLFlow$clone()`

Method `new()`: Creates a new instance of this R6 class.

Usage:

```
OMLFlow$new(id, cache = cache_default(), test_server = test_server_default())
```

Arguments:

`id` (`integer(1)`)

OpenML id for the object.

`cache` (`logical(1) | character(1)`)

See field `cache` for an explanation of possible values. Defaults to value of option `"mlr3oml.cache"`, or `FALSE` if not set.

`test_server` (`character(1)`)

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>). Defaults to value of option `"mlr3oml.test_server"`, or `FALSE` if not set.

Method `print()`: Prints the object.

Usage:

```
OMLFlow$print()
```

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
OMLFlow$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

References

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). “OpenML.” *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
library("mlr3")
# mlr3 flow:
flow = OMLFlow$new(id = 19103)
# using sugar
flow = oml_flow(id = 19103)
learner = as_learner(flow, "classif")
# python flow
python_flow = OMLFlow$new(19090)
# conversion to pseudo Learner
plearner = as_learner(python_flow, "classif")
```

oml_object_r6

Abstract Base Class for OpenML objects.

Description

All OML Objects inherit from this class. Don't use his class directly.

Active bindings

```
desc (list())
  Description of OpenML object.
cache_dir (logical(1) | character(1))
  Stores the location of the cache for objects retrieved from OpenML. If set to FALSE, caching
  is disabled. Objects from the test server are stored in the subdirectory 'test', those from the
  public server are stored in the subdirectory 'public'.
  The package qs is required for caching.
id (integer(1))
  OpenML data id.
server (character(1))
  The server for this object.
man (character(1))
  The manual entry.
name (character(1))
  The name of the object.
type (character())
  The type of OpenML object (e.g. task, run, ...).
test_server (logical(1))
  Whether the object is using the test server.
```


Methods**Public methods:**

- [OMLObject\\$new\(\)](#)
- [OMLObject\\$help\(\)](#)
- [OMLObject\\$clone\(\)](#)

Method `new()`: Creates a new instance of this [R6](#) class.

Usage:

```
OMLObject$new(
  id,
  cache = cache_default(),
  test_server = test_server_default(),
  type
)
```

Arguments:

`id` (`integer(1)`)

OpenML id for the object.

`cache` (`logical(1)` | `character(1)`)

See field `cache` for an explanation of possible values. Defaults to value of option `"mlr3oml.cache"`, or `FALSE` if not set.

`test_server` (`character(1)`)

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>). Defaults to value of option `"mlr3oml.test_server"`, or `FALSE` if not set.

`type` (`character()`)

The type of OpenML object (e.g. `run`, `task`, ...).

Method `help()`: Opens the corresponding help page referenced by field `$man`.

Usage:

```
OMLObject$help()
```

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
OMLObject$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

oml_run_r6

Interface to OpenML Runs

Description

This is the class for OpenML **Runs**, which are conceptually similar to `mlr3::ResampleResults`. This object can also be constructed using the sugar function `oml_run()`.

OpenML Integration

- A `OMLTask` is returned by accessing the active field `$task`.
- A `OMLData` is returned by accessing the active field `$data` (short for `$task$data`)
- A `OMLFlow` is returned by accessing the active field `$flow`.
- The raw predictions are returned by accessing the active field `$prediction`.

mlr3 Integration

- A `mlr3::ResampleResult` is returned when calling `mlr3::as_resample_result()`.
- A `mlr3::Task` is returned when calling `mlr3::as_task()`.
- A `mlr3::DataBackend` is returned when calling `mlr3::as_data_backend()`.
- A instantiated `mlr3::Resampling` is returned when calling `mlr3::as_resampling()`.

Super class

```
mlr3oml::OMLObject -> OMLRun
```

Active bindings

```
flow_id (integer(1))
  The id of the flow.
flow (OMLFlow)
  The OpenML Flow.
tags (character())
  Returns all tags of the object.
parquet (logical(1))
  Whether to use parquet.
task_id (character(1))
  The id of the task solved by this run.
task (OMLTask)
  The task solved by this run.
data_id (integer(1))
  The id of the dataset.
data (OMLData)
  The data used in this run.
```

`task_type` (character())
The task type.

`parameter_setting` data.table()
The parameter setting for this run.

`prediction` (data.table())
The raw predictions of the run as returned by OpenML, not in standard mlr3 format. Formatted predictions are accessible after converting to a [mlr3::ResampleResult](#) via `as_resample_result()`.

`evaluation` (data.table())
The evaluations calculated by the OpenML server.

Methods

Public methods:

- [OMLRun\\$new\(\)](#)
- [OMLRun\\$print\(\)](#)
- [OMLRun\\$clone\(\)](#)

Method `new()`: Creates a new instance of this R6 class.

Usage:

```
OMLRun$new(
  id,
  cache = cache_default(),
  parquet = parquet_default(),
  test_server = test_server_default()
)
```

Arguments:

`id` (integer(1))

OpenML id for the object.

`cache` (logical(1) | character(1))

See field `cache` for an explanation of possible values. Defaults to value of option `"mlr3oml.cache"`, or FALSE if not set.

`parquet` (logical(1))

Whether to use parquet instead of arff. If parquet is not available, it will fall back to arff.

Defaults to value of option `"mlr3oml.parquet"` or FALSE if not set.

`test_server` (character(1))

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>).

Defaults to value of option `"mlr3oml.test_server"`, or FALSE if not set.

Method `print()`: Prints the object.

Usage:

```
OMLRun$print()
```

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
OMLRun$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

References

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). “OpenML.” *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
library("mlr3")
orun = OMLRun$new(id = 10587724)
# sugar
orun = oml_run(id = 10587724)
print(orun)
print(orun$task) # OMLTask
print(orun$data) # OMLData
print(orun$flow) # OMLFlow
print(orun$prediction)
as_task(orun)
as_resampling(orun)
as_data_backend(orun)
rr = as_resample_result(orun)
rr$score(msr("classif.ce"))
```

oml_sugar

Syntactic Sugar for Object Construction

Description

Functions to create OpenML objects. The following functions are available:

- `oml_data()` - creates an instance of the R6 class [OMLData](#).
- `oml_task()` - creates an instance of the R6 class [OMLTask](#).
- `oml_flow()` - creates an instance of the R6 class [OMLFlow](#).
- `oml_run()` - creates an instance of the R6 class [OMLRun](#).
- `oml_collection()` - creates an instance of the R6 class [OMLCollection](#).

Usage

```
oml_data(  
  id,  
  cache = cache_default(),  
  parquet = parquet_default(),  
  test_server = test_server_default()  
)  
  
oml_task(  
  id,
```

```

    cache = cache_default(),
    parquet = parquet_default(),
    test_server = test_server_default()
  )

oml_flow(id, cache = cache_default(), test_server = test_server_default())

oml_run(
  id,
  cache = cache_default(),
  parquet = parquet_default(),
  test_server = test_server_default()
)

oml_collection(
  id,
  cache = cache_default(),
  parquet = parquet_default(),
  test_server = test_server_default()
)

```

Arguments

id	(integer(1)) OpenML id for the object.
cache	(logical(1) character(1)) See field cache for an explanation of possible values. Defaults to value of option "mlr3oml.cache", or FALSE if not set.
parquet	(logical(1)) Whether to use parquet instead of arff. If parquet is not available, it will fall back to arff. Defaults to value of option "mlr3oml.parquet" or FALSE if not set.
test_server	(character(1)) Whether to use the OpenML test server (https://test.openml.org/) or public server (https://www.openml.org/). Defaults to value of option "mlr3oml.test_server", or FALSE if not set.

Value

(OMLObject)

Description

This is the class for tasks served on **OpenML**. It consists of a dataset and other meta-information such as the target variable for supervised problems. This object can also be constructed using the sugar function `oml_task()`.

mlr3 Integration

- Obtain a `mlr3::Task` by calling `as_task()`.
- Obtain a `mlr3::Resampling` by calling `as_resampling()`.

Super class

`mlr3oml::OMLObject` -> `OMLTask`

Active bindings

`estimation_procedure` (`list()`)
The estimation procedure, returns NULL if none is available.

`task_splits` (`data.table()`)
A `data.table` containing the splits as provided by OpenML.

`tags` (`character()`)
Returns all tags of the object.

`parquet` (`logical(1)`)
Whether to use parquet.

`name` (`character(1)`)
Name of the task, extracted from the task description.

`task_type` (`character(1)`)
The OpenML task type.

`data_id` (`integer()`)
Data id, extracted from the task description.

`data` (`OMLData`)
Access to the underlying OpenML data set via a `OMLData` object.

`nrow` (`integer()`)
Number of rows, extracted from the `OMLData` object.

`ncol` (`integer()`)
Number of columns, as extracted from the `OMLData` object.

`target_names` (`character()`)
Name of the targets, as extracted from the OpenML task description.

`feature_names` (`character()`)
Name of the features (without targets of this `OMLTask`).

`data_name` (`character()`)
Name of the dataset (inferred from the task name).

Methods

Public methods:

- [OMLTask\\$new\(\)](#)
- [OMLTask\\$print\(\)](#)
- [OMLTask\\$clone\(\)](#)

Method `new()`: Creates a new instance of this [R6](#) class.

Usage:

```
OMLTask$new(
  id,
  cache = cache_default(),
  parquet = parquet_default(),
  test_server = test_server_default()
)
```

Arguments:

`id` (`integer(1)`)

OpenML id for the object.

`cache` (`logical(1)` | `character(1)`)

See field `cache` for an explanation of possible values. Defaults to value of option `"mlr3oml.cache"`, or `FALSE` if not set.

`parquet` (`logical(1)`)

Whether to use `parquet` instead of `arff`. If `parquet` is not available, it will fall back to `arff`.

Defaults to value of option `"mlr3oml.parquet"` or `FALSE` if not set.

`test_server` (`character(1)`)

Whether to use the OpenML test server (<https://test.openml.org/>) or public server (<https://www.openml.org/>).

Defaults to value of option `"mlr3oml.test_server"`, or `FALSE` if not set.

Method `print()`: Prints the object. For a more detailed printer, convert to a [mlr3::Task](#) via `$task`.

Usage:

```
OMLTask$print()
```

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
OMLTask$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

References

Vanschoren J, van Rijn JN, Bischl B, Torgo L (2014). "OpenML." *ACM SIGKDD Explorations Newsletter*, **15**(2), 49–60. doi: [10.1145/2641190.2641198](https://doi.org/10.1145/2641190.2641198).

Examples

```
library("mlr3")
# Get a task from OpenML:
otask = OMLTask$new(id = 31)
# using sugar
otask = oml_task(id = 31)
otask$data
otask$target_names
otask$feature_names

# convert to mlr3 Task:
task = as_task(otask)

# get a task via tsk():
tsk("oml", task_id = 31L)
```

read_arff

Read ARFF files

Description

Parses a file located at path and returns a `data.table()`.

Limitations:

- Only works for dense files, no support for sparse data. Use **RWeka** instead.
- Dates (even if there is no time component) are read in as `POSIXct`.
- The date-format from the ARFF specification is currently ignored. Instead, we rely on the auto-detection of `data.table`'s `fread()`.

Usage

```
read_arff(path)
```

Arguments

path (character(1))
Path or URI of the ARFF file, passed to `file()`.

Value

`(data.table())`.

write_arff	<i>Write ARFF files</i>
------------	-------------------------

Description

Writes a `data.frame()` to an ARFF file.

Limitations:

- Logicals are written as categorical features.
- [POSIXct](#) columns are converted to UTC.

Usage

```
write_arff(data, path, relation = deparse(substitute(data)))
```

Arguments

data	(<code>data.frame()</code>) Data to write.
path	(<code>character(1)</code>) Path or URI of the ARFF file, passed to <code>file()</code> .
relation	(<code>character(1)</code>) Relation (name) of the data set.

Index

as_learner.OMLFlow, 3

base::make.names(), 11

benchmark_grid_oml, 4

data.table(), 4, 24

data.table::data.table(), 12

data.table::fread(), 2, 11

file(), 24, 25

fread(), 24

list_oml, 4

list_oml_data(list_oml), 4

list_oml_evaluations(list_oml), 4

list_oml_flows(list_oml), 4

list_oml_measures(list_oml), 4

list_oml_runs(list_oml), 4

list_oml_setups(list_oml), 4

list_oml_tasks(list_oml), 4

mlr3::as_benchmark_result, 2, 8

mlr3::as_data_backend, 2

mlr3::as_data_backend(), 11, 18

mlr3::as_learner, 2

mlr3::as_learner(), 15

mlr3::as_learners, 8

mlr3::as_resample_result, 2

mlr3::as_resample_result(), 18

mlr3::as_resampling, 2

mlr3::as_resampling(), 18

mlr3::as_resamplings, 8

mlr3::as_task, 2

mlr3::as_task(), 11, 18

mlr3::as_tasks, 8

mlr3::BenchmarkResult, 8

mlr3::DataBackend, 11, 18

mlr3::DataBackendDataTable, 11

mlr3::Learner, 4, 8, 15

mlr3::mlr_resamplings, 2

mlr3::mlr_tasks, 2

mlr3::ResampleResult, 3, 18, 19

mlr3::Resampling, 2, 4, 8, 18, 22

mlr3::Task, 2, 4, 8, 11, 13, 18, 22, 23

mlr3db::DataBackendDuckDB, 11

mlr3oml (mlr3oml-package), 2

mlr3oml-package, 2

mlr3oml::OMLObject, 8, 12, 15, 18, 22

oml_collection(oml_sugar), 20

oml_collection(), 8

oml_collection_r6, 8

oml_data(oml_sugar), 20

oml_data(), 11

oml_data_r6, 11

oml_flow(oml_sugar), 20

oml_flow(), 14

oml_flow_r6, 14

oml_object_r6, 16

oml_run(oml_sugar), 20

oml_run(), 18

oml_run_r6, 18

oml_sugar, 20

oml_task(oml_sugar), 20

oml_task(), 22

oml_task_r6, 21

OMLCollection, 20

OMLCollection(oml_collection_r6), 8

OMLData, 18, 20, 22

OMLData(oml_data_r6), 11

OMLFlow, 3, 18, 20

OMLFlow(oml_flow_r6), 14

OMLObject, 21

OMLObject(oml_object_r6), 16

OMLRun, 20

OMLRun(oml_run_r6), 18

OMLTask, 4, 18, 20, 22

OMLTask(oml_task_r6), 21

POSIXct, 24, 25

R6, [9](#), [13](#), [15](#), [17](#), [19](#), [23](#)

R_user_dir(), [2](#)

read_arff, [24](#)

RWeka::read.arff(), [11](#)

write_arff, [25](#)