
nesta Documentation

nesta

Nov 03, 2020

Contents:

1	Packages	3
1.1	Code and scripts	3
1.1.1	Meetup	3
1.1.2	Health data	7
1.1.3	NLP Utils	8
1.1.4	Geo Utils	8
1.1.5	Format Utils	9
1.1.6	Decorators	10
1.2	Code auditing	11
2	Production	13
2.1	How to put code into production at nesta	13
2.2	Code and scripts	14
2.2.1	Routines	14
2.2.2	Batchables	22
2.2.3	ORMs	27
2.2.4	Ontologies and schemas	30
2.2.5	Elasticsearch mappings	31
2.2.6	Luigi Hacks	34
2.2.7	Scripts	38
2.2.8	Elasticsearch	39
2.2.9	Containerised Luigi	40
3	FAQ	43
3.1	Where is the data?	43
3.2	Why don't you use Aurora rather than MySQL?	43
3.3	Where are the production machines?	43
3.4	Where is the latest config?	43
3.5	Where do I start with Elasticsearch?	44
4	Troubleshooting	45
4.1	I'm having problems using the config files!	45
4.2	How do I restart the apache server after downtime?	45
4.3	How do I restart the luigi server after downtime?	45
4.4	How do I perform initial setup to ensure the batchables will run?	45
4.5	How can I send/receive emails from Luigi?	46
4.6	How do I add a new user to the server?	46

5 Packages	47
6 Production	49
6.1 License	49
Python Module Index	51
Index	53

Branch	Docs	Build
Master		
Development		

Welcome to [nesta](#)! This repository houses our fully-audited tools and packages, as well as our in-house production system. If you're reading this on [our GitHub repo](#), you will find complete documentation at our [Read the Docs site](#).

Nesta's collection of tools for meaty tasks. Any processes that go into production come here first, but there are other good reasons for code to end up here.

1.1 Code and scripts

1.1.1 Meetup

NB: The meetup pipeline will not work until [this issue](#) has been resolved.

Data collection of Meetup data. The procedure starts with a single country and [Meetup category](#). All of the groups within the country are discovered, from which all members are subsequently retrieved (no personal information!). In order to build a fuller picture, all other groups to which the members belong are retrieved, which may be in other categories or countries. Finally, all group details are retrieved.

The code should be executed in the following order, which reflects the latter procedure:

- 1) country_groups.py
- 2) groups_members.py
- 3) members_groups.py
- 4) groups_details.py

Each script generates a list of dictionaries which can be ingested by the proceeding script.

Country → Groups

Start with a country (and Meetup category) and end up with Meetup groups.

generate_coords ($x0, y0, x1, y1, n$)

Generate $O(\frac{n^2}{2})$ coordinates in the bounding box $(x0, y0), (x1, y1)$, such that overlapping circles of equal radii (situated at each coordinate) entirely cover the area of the bounding box. The longitude and latitude are treated

as euclidean variables, although the radius (calculated from the smallest side of the bounding box divided by n) is calculated correctly. In order for the circles to fully cover the region, an unjustified factor of 10% is included in the radius. Feel free to do the maths and work out a better strategy for covering a geographical area with circles.

The circles (centred on each X) are staggered as so (single vertical lines or four underscores correspond to a circle radius):

```

____X____X____
|
X____X____X
|
____X____X____

```

This configuration corresponds to $n = 4$.

Parameters

- **x0, y0, x1, y1** (*float*) – Bounding box coordinates (lat/lon)
- **n** (*int*) – The fraction by which to calculate the Meetup API radius parameter, with respect to the smallest side of the country's shape bbox. This will generate $\mathcal{O}(\frac{n^2}{2})$ separate Meetup API radius searches. The total number of searches scales with the ratio of the bbox sides.

Returns The radius and coordinates for the Meetup API request

Return type float, list of tuple

`get_coordinate_data(n)`

Generate the radius and coordinate data (see `generate_coords`) for each shape (country) in the shapefile pointed to by the environmental variable `WORLD_BORDERS`.

Parameters **n** (*int*) – The fraction by which to calculate the Meetup API radius parameter, with respect to the smallest side of the country's shape bbox. This will generate $\mathcal{O}(\frac{n^2}{2})$ separate Meetup API radius searches. The total number of searches scales with the ratio of the bbox sides.

Returns

containing coordinate and radius for each country.

Return type `pd.DataFrame`

`assert_iso2_key(df, iso2)`

`class MeetupCountryGroups(country_code, coords, radius, category, n=10)`

Bases: `object`

Extract all meetup groups for a given country.

country_code

ISO2 code

Type str

params (

obj:'dict'): GET request parameters, including lat/lon.

groups

List of meetup groups in this country, assigned assigned after calling `get_groups`.

Type list of str

get_groups (*lon, lat, offset=0, max_pages=None*)

Recursively get all groups for the given parameters. It is assumed that you will run with the default arguments, since they are set automatically in the recursing procedure.

get_groups_recursive ()

Call `get_groups` for each lat,lon coordinate

Groups → Members

Start with Meetup groups and end up with Meetup members.

get_members (*params*)

Hit the Meetup API for the members of a specified group.

Parameters *params* (dict) – `https://api.meetup.com/members/` parameters

Returns Meetup member IDs

Return type (list of str)

get_all_members (*group_id, group_urlname, max_results, test=False*)

Get all of the Meetup members for a specified group.

Parameters

- **group_id** (*int*) – The Meetup ID of the group.
- **group_urlname** (*str*) – The URL name of the group.
- **max_results** (*int*) – The maximum number of results to return per API query.
- **test** (*bool*) – For testing.

Returns A matchable list of Meetup members

Return type (list of dict)

Members → Groups

Start with Meetup members and end up with Meetup groups.

exception NoMemberFound (*member_id*)

Bases: `Exception`

Exception should no member be found by the Meetup API

get_member_details (*member_id, max_results*)

Hit the Meetup API for details of a specified member

Parameters

- **member_id** (*str*) – A Meetup member ID
- **max_results** (*int*) – The maximum number of results with each API hit

Returns Meetup API response json.

Return type list of dict

get_member_groups (*member_info*)

Extract the groups data from Meetup membership information.

Parameters

- **member_id** (*str*) – A Meetup member ID

- **member_info** (list of dict) – Meetup member API response json.

Returns List of unique member-group combinations

Return type list of dict

Groups → Group details

Start with Meetup groups and end up with Meetup group details.

exception NoGroupFound (group_urlname)

Bases: Exception

Exception should no group be found by the Meetup API

get_group_details (group_urlname, max_results, avoid_exception=True)

Hit the Meetup API for the details of a specified groups. :param group_urlname: A Meetup group urlname :type group_urlname: str :param max_results: Total number of results to return per API request. :type max_results: int

Returns Meetup API response data

Return type (list of dict)

Utils

Common tools between the different data collection points.

get_api_key ()

Get a random API key from those listed in the environmental variable MEETUP_API_KEYS.

save_sample (json_data, filename, k)

Dump a sample of k items from row-oriented JSON data json_data into file with name filename.

flatten_data (list_json_data, keys, **kwargs)

Flatten nested JSON data from a list of JSON objects, by a list of desired keys. Each element in the keys may also be an ordered iterable of keys, such that subsequent keys describe a path through the JSON to desired value. For example in order to extract *key1* and *key3* from:

```
{'key': <some_value>, 'key2' : {'key3': <some_value>}}
```

one would specify keys as:

```
['key1', ('key2', 'key3')]
```

Parameters

- **list_json_data** (json) – Row-orientated JSON data.
- **keys** (list) – Mixed list of either: individual *str* keys for data values
- **are not nested; or sublists of str, as described above.** (*which*) –
- ****kwargs** – Any constants to include in every flattened row of the output.

Returns Flattened row-orientated JSON data.

Return type json

get_members_by_percentile (*engine*, *perc=10*)

Get the number of meetup group members for a given percentile from the database.

Parameters

- **engine** – A SQL alchemy connectable.
- **perc** (*int*) – A percentile to evaluate.

Returns The number of members corresponding to this percentile.

Return type members (float)

get_core_topics (*engine*, *core_categories*, *members_limit*, *perc=99*)

Get the most frequent topics from a selection of meetup categories, from the database.

Parameters

- **engine** – A SQL alchemy connectable.
- **core_categories** (*list*) – A list of category_shortnames.
- **members_limit** (*int*) – Minimum number of members required in a group for it to be considered.
- **perc** (*int*) – A percentile to evaluate the most frequent topics.

Returns The set of most frequent topics.

Return type topics (set)

1.1.2 Health data

Initially for our project with the Robert Woods Johnson Foundation (RWJF), these procedures outline the data collection of health-specific data.

Collect NIH

Extract all of the NIH World RePORTER data via their static data dump. N_TABS outputs are produced in CSV format (concatenated across all years), where N_TABS correspondes to the number of tabs in the main table found at:

https://exporter.nih.gov/ExPORTER_Catalog.aspx

The data is transferred to the Nesta intermediate data bucket.

get_data_urls (*tab_index*)

Get all CSV URLs from the `tab_index`'th tab of the main table found at `:code:`TOP_URL``.

Parameters **tab_index** (*int*) – Tab number (0-indexed) of table to extract CSV URLs from.

Returns Title of the tab in the table. hrefs (list): List of URLs pointing to data CSVs.

Return type title (str)

iterrows (*url*)

Yield rows from the CSV (found at URL `url`) as JSON (well, dict objects).

Parameters **url** (*str*) – The URL at which a zipped-up CSV is found.

Yields dict object, representing one row of the CSV.

Process NIH

Data cleaning and processing procedures for the NIH World Reporter data. Specifically, a lat/lon is generated for each city/country; and the formatting of date fields is unified.

1.1.3 NLP Utils

Standard tools for aiding natural language processing.

Preprocess

Tools for preprocessing text.

tokenize_document (*text*, *remove_stops=False*)

Preprocess a whole raw document. :param text: Raw string of text. :type text: str :param remove_stops: Flag to remove english stopwords :type remove_stops: bool

Returns List of preprocessed and tokenized documents

clean_and_tokenize (*text*, *remove_stops*)

Preprocess a raw string/sentence of text. :param text: Raw string of text. :type text: str :param remove_stops: Flag to remove english stopwords :type remove_stops: bool

Returns Preprocessed tokens.

Return type tokens (list, str)

filter_by_idf (*documents*, *lower_idf_limit*, *upper_idf_limit*)

Remove (from documents) terms which are in a range of IDF values.

Parameters

- **documents** (*list*) – Either a list of str or a list of list of str to be filtered.
- **lower_idf_limit** (*float*) – Lower percentile (between 0 and 100) on which to exclude terms by their IDF.
- **upper_idf_limit** (*float*) – Upper percentile (between 0 and 100) on which to exclude terms by their IDF.

Returns Filtered documents

1.1.4 Geo Utils

Tools for processing of geographical data, such as geocoding.

geocode

Tools for geocoding.

geocode (***request_kwargs*)

Geocoder using the Open Street Map Nominatim API.

If there are multiple results the first one is returned (they are ranked by importance). The API usage policy allows maximum 1 request per second and no multithreading: <https://operations.osmfoundation.org/policies/nominatim/>

Parameters **request_kwargs** (*dict*) – Parameters for OSM API.

Returns JSON from API response.

retry_if_not_value_error (*exception*)

Forces retry to exit if a ValueError is returned. Supplied to the 'retry_on_exception' argument in the retry decorator.

Parameters **exception** (*Exception*) – the raised exception, to check

Returns False if a ValueError, else True

Return type (bool)

geocode_dataframe (*df*)

A wrapper for the geocode function to process a supplied dataframe using the city and country.

Parameters **df** (*dataframe*) – a dataframe containing city and country fields.

Returns a dataframe with a 'coordinates' column appended.

geocode_batch_dataframe (*df*, *city='city'*, *country='country'*, *latitude='latitude'*, *longitude='longitude'*, *query_method='both'*)

Geocodes a dataframe, first by supplying the city and country to the api, if this fails a second attempt is made supplying the combination using the q= method. The supplied dataframe df is returned with additional columns appended, containing the latitude and longitude as floats.

Parameters

- **df** (*pandas.DataFrame*) – input dataframe
- **city** (*str*) – name of the input column containing the city
- **country** (*str*) – name of the input column containing the country
- **latitude** (*str*) – name of the output column containing the latitude
- **longitude** (*str*) – name of the output column containing the longitude
- **query_method** (*int*) – query methods to attempt: 'city_country_only': city and country only 'query_only': q method only 'both': city, country with fallback to q method

Returns original dataframe with lat and lon appended as floats

Return type (*pandas.DataFrame*)

generate_composite_key (*city=None*, *country=None*)

Generates a composite key to use as the primary key for the geographic data.

Parameters

- **city** (*str*) – name of the city
- **country** (*str*) – name of the country

Returns composite key

Return type (*str*)

1.1.5 Format Utils

Tools for formatting data, such as dates.

datetools

Tools for processing dates in data.

extract_year (*date*)

Use search for 4 digits in a row to identify the year and return as YYYY-01-01.

Parameters **date** (*str*) – The full date string.

Returns integer

extract_date (*date*, *date_format*='%Y-%m-%d', *return_date_object*=False)

Determine the date format, convert and return in YYYY-MM-DD format.

Parameters **date** (*str*) – the full date string.

Returns Formatted date string.

1.1.6 Decorators

ratelimit

Apply rate limiting at a threshold per second

ratelimit (*max_per_second*)

Parameters **max_per_second** (*float*) – Number of permitted hits per second

schema_transform

Apply a field name transformation to a data output from the wrapped function, such that specified field names are transformed and unspecified fields are dropped. A valid file would be formatted as shown:

```
{ "tier0_to_tier1":
  { "bad_col": "good_col", "another_bad_col": "another_good_col"
  }
}
```

load_transformer (*filename*)

schema_transform (*filename*)

Parameters **filename** (*str*) – A record-oriented JSON file path mapping field names

Returns Data in the format it was originally passed to the wrapper in, with specified field names transformed and unspecified fields dropped.

schema_transformer (*data*, *, *filename*, *ignore*=[])

Function version of the schema_transformer wrapper. :param data: the data requiring the schama transformation :type data: dataframe OR list of dicts :param filename: the path to the schema json file :type filename: str :param ignore: optional list of fields, eg ids or keys which shouldn't be dropped :type ignore: list

Returns supplied data with schema applied

1.2 Code auditing

Packages are only accepted if they satisfy our internal auditing procedure:

- **Common sense requirements:**
 - **Either:**
 - * The code produces at least one data or model output; **or**
 - * The code provides a service which abstracts away significant complexity.
 - There is one unit test for each function or method, which lasts no longer than about 1 minute.
 - Each data or model output is produced from a single function or method, as described in the `__main__` of a specified file.
 - Can the nearest programmer (or equivalent) checkout and execute your tests from scratch?
 - Will the code be used to perform non-project specific tasks?
 - Does the process perform a logical task or fulfil a logical purpose?
- **If the code requires productionising, it satisfies one of the following conditions:**
 - a) There is a non-trivial pipeline, which would benefit from formal productionising.
 - b) A procedure is foreseen to be reperformed for new contexts with atomic differences in run conditions.
 - c) The output is a service which requires a pipeline.
 - d) The process is a regular / longitudinal data collection.
- **Basic PEP8 and style requirements:**
 - Docstrings for every exposable class, method and function.
 - Usage in a README.rst or in Docstring at the top of the file.
 - CamelCase class names.
 - Underscore separation of all other variable, function and module names.
 - No glaring programming no-nos.
 - Never use `print: opt` for logging instead.
- **Bureaucratic requirements:**
 - A requirements file*.
 - The README file specifies the operating system and python version.

Nesta's production system is based on [Luigi](#) pipelines, and are designed to be entirely run on AWS via the batch service. The main Luigi server runs on a persistent EC2 instance. Beyond the well documented Luigi code, the main features of the nesta production system are:

- `luigihacks.autobatch`, which facilitates a managed `Luigi.Task` which is split, batched and combined in a single step. Currently only synchronous jobs are accepted. Asynchronous jobs (where downstream `Luigi.Task` jobs can be triggered) are a part of a longer term plan.
- `scripts.nesta_prepare_batch` which zips up the batchable with the specified environmental files and ships it to AWS S3.
- `scripts.nesta_docker_build` which builds a specified docker environment and ships it to AWS ECS.

2.1 How to put code into production at nesta

If you're completely new, check out our [training slides](#). In short, the steps you should go through when building production code are to:

1. Audit the package code, required to pass all auditing tests
2. Understand what environment is required
3. Write a Dockerfile and docker launch script for this under `scripts/docker_recipes`
4. Build the Docker environment (run: `docker_build <recipe_name>` from any directory)
5. Build and test the batchable(s)
6. Build and test a Luigi pipeline
7. [...] Need to have steps here which estimate run time cost parameters. Could use `tests.py` to estimate this. [...]
8. Run the full chain

2.2 Code and scripts

2.2.1 Routines

All of our pipelines, implemented as Luigi routines. Some of these pipelines (at least partly) rely on batch computing (via AWS Batch), where the ‘batched’ scripts (*run.py* modules) are described in `core.batchables`. Other than `luigihacks.autobatch`, which is respectively documented, the routine procedure follows the core Luigi documentation.

Examples

Examples of Luigi routines, from which all other nesta production routines can be built. Currently we have examples of routines with S3 and database (MySQL) IO, and routines which are entirely batched.

We’d recommend reading Spotify’s Luigi documentation, and also checking the *Luigi Hacks* documentation which contains modified Luigi modules which (who knows) one day we will suggest as pull requests.

S3 Example

An example of building a pipeline with S3 Targets

```
class InputData (*args, **kwargs)
    Bases: luigi.task.ExternalTask

    Dummy task acting as the single input data source

    output ()
        Points to the S3 Target

class SomeTask (*args, **kwargs)
    Bases: luigi.task.Task

    An intermediate task which increments the age of the muppets by 1 year.

    Parameters date (datetime) – Date used to label the outputs

    date = <luigi.parameter.DateParameter object>

    requires ()
        Gets the input data (json containing muppet name and age)

    output ()
        Points to the S3 Target

    run ()
        Increments the muppets’ ages by 1

class FinalTask (*args, **kwargs)
    Bases: luigi.task.Task

    The root task, which adds the surname ‘Muppet’ to the names of the muppets.

    Parameters date (datetime) – Date used to label the outputs

    date = <luigi.parameter.DateParameter object>

    requires ()
        Get data from the intermediate task
```

```

output ()
    Points to the S3 Target

run ()
    Appends 'Muppet' the muppets' names
  
```

Database example

An example of building a pipeline with database Targets

```

class InputData (*args, **kwargs)
    Bases: luigi.task.Task

    Dummy task acting as the single input data source.

    Parameters
    • date (datetime) – Date used to label the outputs
    • db_config – (dict) The input database configuration

    date = <luigi.parameter.DateParameter object>
    db_config = <luigi.parameter.DictParameter object>

    output ()
        Points to the input database target

    run ()
        Example of marking the update table

class SomeTask (*args, **kwargs)
    Bases: luigi.task.Task

    Task which increments the age of the muppets, by first selecting muppets with an age less than max_age.

    Parameters
    • date (datetime) – Date used to label the outputs
    • max_age (int) – Maximum age of muppets to select from the database
    • in_db_config – (dict) The input database configuration
    • out_db_config – (dict) The output database configuration

    date = <luigi.parameter.DateParameter object>
    max_age = <luigi.parameter.IntParameter object>
    in_db_config = <luigi.parameter.DictParameter object>
    out_db_config = <luigi.parameter.DictParameter object>

    requires ()
        Gets the input database engine

    output ()
        Points to the output database engine

    run ()
        Increments the muppets' ages by 1
  
```

```
class RootTask (*args, **kwargs)
```

```
    Bases: luigi.task.WrapperTask
```

A dummy root task, which collects the database configurations and executes the central task.

Parameters `date` (*datetime*) – Date used to label the outputs

```
date = <luigi.parameter.DateParameter object>
```

```
requires ()
```

Collects the database configurations and executes the central task.

arXiv data (technical research)

Data collection and processing pipeline for arXiv data, principally for the [arXlive](#) platform. This pipeline orchestrates the collection of arXiv data, enrichment (via MAG and GRID), topic modelling, and novelty (lolvelty) measurement.

Collection task

Luigi routine to collect new data from the arXiv api and load it to MySQL.

```
class CollectNewTask (*args, **kwargs)
```

```
    Bases: luigi.task.Task
```

Collect new data from the arXiv api and dump the data in the MySQL server.

Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **_routine_id** (*str*) – String used to label the AWS task
- **db_config_env** (*str*) – environmental variable pointing to the db config file
- **db_config_path** (*str*) – The output database configuration
- **insert_batch_size** (*int*) – number of records to insert into the database at once
- **articles_from_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format

```
date = <luigi.parameter.DateParameter object>
```

```
test = <luigi.parameter.BoolParameter object>
```

```
db_config_env = <luigi.parameter.Parameter object>
```

```
db_config_path = <luigi.parameter.Parameter object>
```

```
insert_batch_size = <luigi.parameter.IntParameter object>
```

```
articles_from_date = <luigi.parameter.Parameter object>
```

```
output ()
```

Points to the output database engine

```
run ()
```

The task run method, to be overridden in a subclass.

See Task.run

Date task

Luigi wrapper to identify the date since the last iterative data collection

```
class DateTask(*args, **kwargs)
```

```
    Bases: luigi.task.WrapperTask
```

Collect new data from the arXiv api and dump the data in the MySQL server.

Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **_routine_id** (*str*) – String used to label the AWS task
- **db_config_env** (*str*) – environmental variable pointing to the db config file
- **db_config_path** (*str*) – The output database configuration
- **insert_batch_size** (*int*) – number of records to insert into the database at once
- **articles_from_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format

```
date = <luigi.parameter.DateParameter object>
```

```
test = <luigi.parameter.BoolParameter object>
```

```
db_config_path = <luigi.parameter.Parameter object>
```

```
db_config_env = <luigi.parameter.Parameter object>
```

```
insert_batch_size = <luigi.parameter.IntParameter object>
```

```
articles_from_date = <luigi.parameter.Parameter object>
```

```
requires()
```

Collects the last date of successful update from the database and launches the iterative data collection task.

arXiv enriched with MAG (API)

Luigi routine to query the Microsoft Academic Graph for additional data and append it to the exiting data in the database.

```
class QueryMagTask(*args, **kwargs)
```

```
    Bases: luigi.task.Task
```

Query the MAG for additional data to append to the arxiv articles, primarily the fields of study.

Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **_routine_id** (*str*) – String used to label the AWS task
- **db_config_env** (*str*) – environmental variable pointing to the db config file
- **db_config_path** (*str*) – The output database configuration
- **mag_config_path** (*str*) – Microsoft Academic Graph Api key configuration path
- **insert_batch_size** (*int*) – number of records to insert into the database at once (not used in this task but passed down to others)

- **articles_from_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```

date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>

```

output ()
Points to the output database engine

requires ()
The Tasks that this Task depends on.

A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

See Task.requires

run ()
The task run method, to be overridden in a subclass.
See Task.run

arXiv enriched with MAG (SPARQL)

Luigi routine to query the Microsoft Academic Graph for additional data and append it to the exiting data in the database. This is to collect information which is difficult to retrieve via the MAG API.

```

class MagSparqlTask (*args, **kwargs)
    Bases: luigi.task.Task

```

Query the MAG for additional data to append to the arxiv articles, primarily the fields of study.

Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **_routine_id** (*str*) – String used to label the AWS task
- **db_config_env** (*str*) – environmental variable pointing to the db config file
- **db_config_path** (*str*) – The output database configuration
- **mag_config_path** (*str*) – Microsoft Academic Graph Api key configuration path
- **insert_batch_size** (*int*) – number of records to insert into the database at once (not used in this task but passed down to others)
- **articles_from_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```

date = <luigi.parameter.DateParameter object>

```

```

test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>

output ()
    Points to the output database engine

requires ()
    The Tasks that this Task depends on.

    A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

    See Task.requires

run ()
    The task run method, to be overridden in a subclass.

    See Task.run

```

arXiv enriched with GRID

Luigi routine to lookup arXiv author's institutes via the GRID data, in order to "geocode" arXiv articles. The matching of institute name to GRID data is done via smart(ish) fuzzy matching, which then gives a confidence score per match.

```
class GridTask (*args, **kwargs)
```

Bases: `luigi.task.Task`

Join arxiv articles with GRID data for institute addresses and geocoding.

Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **_routine_id** (*str*) – String used to label the AWS task
- **db_config_env** (*str*) – environmental variable pointing to the db config file
- **db_config_path** (*str*) – The output database configuration
- **mag_config_path** (*str*) – Microsoft Academic Graph Api key configuration path
- **insert_batch_size** (*int*) – number of records to insert into the database at once (not used in this task but passed down to others)
- **articles_from_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```

date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>

```

`insert_batch_size = <luigi.parameter.IntParameter object>`

`articles_from_date = <luigi.parameter.Parameter object>`

`output ()`

Points to the output database engine

`requires ()`

The Tasks that this Task depends on.

A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

See Task.requires

`run ()`

The task run method, to be overridden in a subclass.

See Task.run

`class GridRootTask (*args, **kwargs)`

Bases: `luigi.task.WrapperTask`

`date = <luigi.parameter.DateParameter object>`

`db_config_path = <luigi.parameter.Parameter object>`

`production = <luigi.parameter.BoolParameter object>`

`drop_and_recreate = <luigi.parameter.BoolParameter object>`

`articles_from_date = <luigi.parameter.Parameter object>`

`insert_batch_size = <luigi.parameter.IntParameter object>`

`debug = <luigi.parameter.BoolParameter object>`

`requires ()`

Collects the database configurations and executes the central task.

CORDIS (EU funded research)

Generic pipeline (i.e. not project specific) to collect all CORDIS data, discovering all entities by crawling an unofficial API.

Crunchbase (private sector companies)

NB: The Crunchbase pipeline may not work until [this issue](#) has been resolved.

Data collection and processing pipeline of Crunchbase data, principally for the [healthMosaic](#) platform.

EURITO (piping data to Elasticsearch)

Pipeline specific to EURITO for piping existing data to Elasticsearch. A recent "EU" cut of patstat data is transferred from the "main" patstat database, to Nesta's central database.

Gateway to Research (UK publicly funded research)

Generic pipeline (i.e. not project specific) to collect all GtR data, discovering all entities by crawling the official API. The routine then geocodes and loads data to MYSQL.

NiH (health research)

Data collection and processing pipeline of NiH data, principally for the [healthMosaic](#) platform.

Meetup (social networking data)

NB: The Meetup pipeline will not work until [this issue](#) has been resolved.

Data collection and processing pipeline of Meetup data, principally for the [healthMosaic](#) platform.

Topic discovery

Task to automatically discover relevant topics from meetup data, defined as the most frequently occurring from a set of categories.

```
class TopicDiscoveryTask (*args, **kwargs)
    Bases: luigi.task.Task
```

Task to automatically discover relevant topics from meetup data, defined as the most frequently occurring from a set of categories.

Parameters

- **db_config_env** (*str*) – Environmental variable pointing to the path of the DB config.
- **routine_id** (*str*) – The routine UID.
- **core_categories** (*list*) – A list of category_shortnames from which to identify topics.
- **members_perc** (*int*) – A percentile to evaluate the minimum number of members.
- **topic_perc** (*int*) – A percentile to evaluate the most frequent topics.
- **test** (*bool*) – Test mode.

```
db_config_env = <luigi.parameter.Parameter object>
routine_id = <luigi.parameter.Parameter object>
core_categories = <luigi.parameter.ListParameter object>
members_perc = <luigi.parameter.IntParameter object>
topic_perc = <luigi.parameter.IntParameter object>
test = <luigi.parameter.BoolParameter object>

output ()
    Points to the S3 Target

run ()
    Extract the topics of interest
```

2.2.2 Batchables

Packets of code to be batched by `core.routines` routines. Each packet should sit in it's own directory, with a file called `run.py`, containing a 'main' function called `run()` which will be executed on the AWS batch system.

Each `run.py` should expect an environment parameter called `BATCHPAR_outfile` which should provide information on the output location. Other input parameters should be prefixed with `BATCHPAR_`, as set in `core.routines` routine.

Data / project specific batchables

Example

There are two batchable examples listed here. The first is a module which will be run if you try executing the `batch_example` luigi routine. The second is purely meant as a template, if you are learning the design pattern for nesta's luigi batchables.

`run.py (batch_example)`

The batchable for the `routines.examples.batch_example`, which simply increments a muppet's age by one unit.

run()

Gets the name and age of the muppet, and increments the age. The result is transferred to S3.

`run.py (template_batchable)`

This is a pretty generic example of how your `run.py` might look. It reads and writes from a table, and hits the S3 "checkpoint" at the end.

run()

arXiv data (technical research)

CORDIS (EU-funded research)

`run.py (cordis_api)`

Transfer data on organisations, projects and outputs from the Cordis API on a project-by-project basis.

extract_core_orgs (*orgs*, *project_rcn*)

Seperate a project-organisation (which) is likely to be a department, with a non-unique address.

Parameters

- **orgs** (*list*) – List of organisations to process (NB: this will be modified)
- **project_rcn** (*str*) – The record number of this project

Returns The unique 'parent' organisations.

Return type `core_orgs` (*list*)

prepare_data (*items, rcn*)

Append the project code ('RCN') to each "row" (dict) of data (list)

split_links (*items, project_rcn*)

Generate link table items for each item (dict) in items (list) for the project

run ()

Crunchbase data (private companies)

NB: The Crunchbase pipeline may not work until [this issue](#) has been resolved.

Batchables for the collection and processing of Crunchbase data. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The data is collected from proprietary data dumps, parsed into MySQL (tier 0) and then piped into Elasticsearch (tier 1), post-processing.

run.py (crunchbase_collect)

Collect Crunchbase data from the proprietary data dump and pipe into the MySQL database.

run ()

EURITO

Batchables for processing data (which has already been collected elsewhere within this codebase) for the EURITO project. All of these batchables pipe the data into an Elasticsearch database, which is then cloned by EURITO.

run.py (patstat_eu)

Transfer pre-collected PATSTAT data from MySQL to Elasticsearch. Only EU patents since the year 2000 are considered. The patents are grouped by patent families.

select_text (*objs, lang_field, text_field*)

metadata (*orm, session, appln_ids, field_selector=None*)

run ()

GtR (UK publicly funded research)

Batchable tools for collecting and processing GtR data. As documented under packages and routines, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The data is collected by traversing the graph exposed by the GtR API, and is parsed into MySQL (tier 0). There is a further module for directly generating document embeddings of GtR project descriptions, which can be used for finding topics.

run.py (collect_gtr)

Starting from GtR projects, iteratively and recursively discover all GtR entities by crawling the API.

run ()

run.py (embed_topics)

Document embedding of GtR data. Would be better if this was generalized (i.e. not GtR specific), and migrated to batchables.nlp [see <https://github.com/nestauk/nesta/issues/203>]

run ()

NiH data (health research)

Batchables for the collection and processing of NiH data. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The data is collected from official data dumps, parsed into MySQL (tier 0) and then piped into Elasticsearch (tier 1), post-processing.

run.py (nih_process_data)

Geocode NiH data (from MySQL) and pipe into Elasticsearch.

run ()

run.py (nih_abstract_mesh_data)

Retrieve NiH abstracts from MySQL, assign pre-calculated MeSH terms for each abstract, and pipe data into Elasticsearch. Exact abstract duplicates are removed at this stage.

clean_abstract (*abstract*)

Removes multiple spaces, tabs and newlines.

Parameters **abstract** (*str*) – text to be cleaned

Returns (*str*): cleaned text

run ()

run.py (nih_dedupe)

Deduplicate NiH articles based on similarity scores using Elasticsearch's document similarity API. Similarity is calculated based on the description of the project, the project abstract and the title of the project. Funding information is aggregated (summed) across all deduplicated articles, for the total and annuals funds.

get_value (*obj*, *key*)

Retrieve a value by key if exists, else return None.

extract_yearly_funds (*src*)

Extract yearly funds

run ()

Meetup (social networking / knowledge exchange)

NB: The meetup pipeline will not work until [this issue](#) has been resolved.

Batchables for the Meetup data collection pipeline. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The `topic_tag_elasticsearch` module is responsible for piping data to elasticsearch, as well as apply topic tags and filtering small groups out of the data.

run.py (country_groups)

Batchable for expanding from countries to groups

run ()

run.py (groups_members)

Batchable for expanding group members

run ()

run.py (members_groups)

Batchable for expanding from members to groups.

run ()

run.py (group_details)

Batchable for expanding group details

run ()

run.py (topic_tag_elasticsearch)

Batchable for piping data to Elasticsearch, whilst implementing topic tags, and filtering groups with too few members (given by the 10th percentile of group size, to avoid "junk" groups).

run ()

General-purpose batchables

Bulk geocoding

run.py (batch_geocode)

Geocode any row-delimited json data, with columns corresponding to a city/town/etc and country.

run ()

Natural Language Processing

Batchable utilities for NLP. Note that modules prefixed with [AutoML] are designed to be launched by AutoMLTask, and those with the addition * (i.e. [AutoML*]) are the designed to be the final task in an AutoMLTask chain (i.e. they provide a 'loss').

[AutoML*] run.py (corex_topic_model)

Generate topics based on the CorEx algorithm. Loss is calculated from the total correlation explained.

run ()

[AutoML] run.py (ngrammer)

Find and replace ngrams in a body of text, based on Wiktionary N-Grams. Whilst at it, the ngrammer also tokenizes and removes stop words (unless they occur within an n-gram)

run ()

[AutoML] run.py (tfidf)

Applies TFIDF cuts to a dataset via environmental variables lower_tfidf_percentile and upper_tfidf_percentile.

chunker (*_transformed, n_chunks*)
Yield chunks from a numpy array.

Parameters

- **_transformed** (*np.array*) – Array to split into chunks.
- **n_chunks** (*int*) – Number of chunks to split the array into.

Yields chunk (*np.array*)

run ()

[AutoML] vectorizer (run.py)

Vectorizes (counts or binary) text data, and applies basic filtering of extreme term/document frequencies.

term_counts (*dct, row, binary=False*)
Convert a single single document to term counts via a gensim dictionary.

Parameters

- **dct** (*Dictionary*) – Gensim dictionary.
- **row** (*str*) – A document.
- **binary** (*bool*) – Binary rather than total count?

Returns dict of term id (from the Dictionary) to term count.

optional (*name, default*)

Defines optional env fields with default values

merge_lists (*list_of_lists*)

Join a lists of lists into a single list. Returns an empty list if the input is not a list, which is expected to happen (from the ngrammer) if no long text was found

run ()

Novelty

Batchables for calculating measures of “novelty”.

run.py (lolvelty)

Calculates the “lolvelty” novelty score to documents in Elasticsearch, on a document-by-document basis. Note that this is a slow procedure, and the bounds of document “lolvelty” can’t be known a priori.

run ()

2.2.3 ORMs

SQLAlchemy ORMs for the routines, which allows easy integration of testing (including automatic setup of test databases and tables).

Meetup

```
class Group (**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    id
    urlname
    category_name
    category_shortcode
    city
    country
    created
    description
    lat
    lon
    members
    name
    topics
    category_id
```

`country_name`

`timestamp`

class `GroupMember` (***kwargs*)

Bases: `sqlalchemy.ext.declarative.api.Base`

Note: no foreign key constraint, since unknown groups will be found in the member expansion phase

`group_id`

`group_urlname`

`member_id`

NIH schema

The schema for the World RePORTER data.

class `Projects` (***kwargs*)

Bases: `sqlalchemy.ext.declarative.api.Base`

`application_id`

`activity`

`administering_ic`

`application_type`

`arra_funded`

`award_notice_date`

`budget_start`

`budget_end`

`cfda_code`

`core_project_num`

`ed_inst_type`

`foa_number`

`full_project_num`

`funding_ics`

`funding_mechanism`

`fy`

`ic_name`

`org_city`

`org_country`

`org_dept`

`org_district`

`org_duns`

`org_fips`

`org_ipf_code`


```

    org_name
    org_state
    org_zipcode
    phr
    pi_ids
    pi_names
    program_officer_name
    project_start
    project_end
    project_terms
    project_title
    serial_number
    study_section
    study_section_name
    suffix
    support_year
    direct_cost_amt
    indirect_cost_amt
    total_cost
    subproject_id
    total_cost_sub_project
    nih_spending_cats

class Abstracts(**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    application_id
    abstract_text

class Publications(**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    pmid
    author_name
    affiliation
    author_list
    country
    issn
    journal_issue
    journal_title
    journal_title_abbr

```

```

    journal_volume
    lang
    page_number
    pub_date
    pub_title
    pub_year
    pmc_id
class Patents (**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    patent_id
    patent_title
    project_id
    patent_org_name
class LinkTables (**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    pmid
    project_number
class ClinicalStudies (**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    clinicaltrials_gov_id
    core_project_number
    study
    study_status

```

2.2.4 Ontologies and schemas

Tier 0

Raw data collections (“tier 0”) in the production system do not adhere to a fixed schema or ontology, but instead have a schema which is very close to the raw data. Modifications to field names tend to be quite basic, such as lowercase and removal of whitespace in favour of a single underscore.

Tier 1

Processed data (“tier 1”) is intended for public consumption, using a common ontology. The convention we use is as follows:

- Field names are composed of up to three terms: a `firstName`, `middleName` and `lastName`
- Each term (e.g. `firstName`) is written in `lowerCamelCase`.
- `firstName` terms correspond to a restricted set of basic quantities.

- `middleName` terms correspond to a restricted set of modifiers (e.g. adjectives) which add nuance to the `firstName` term. Note, the special `middleName` term `of` is reserved as the default value in case no `middleName` is specified.
- `lastName` terms correspond to a restricted set of entity types.

Valid examples are `date_start_project` and `title_of_project`.

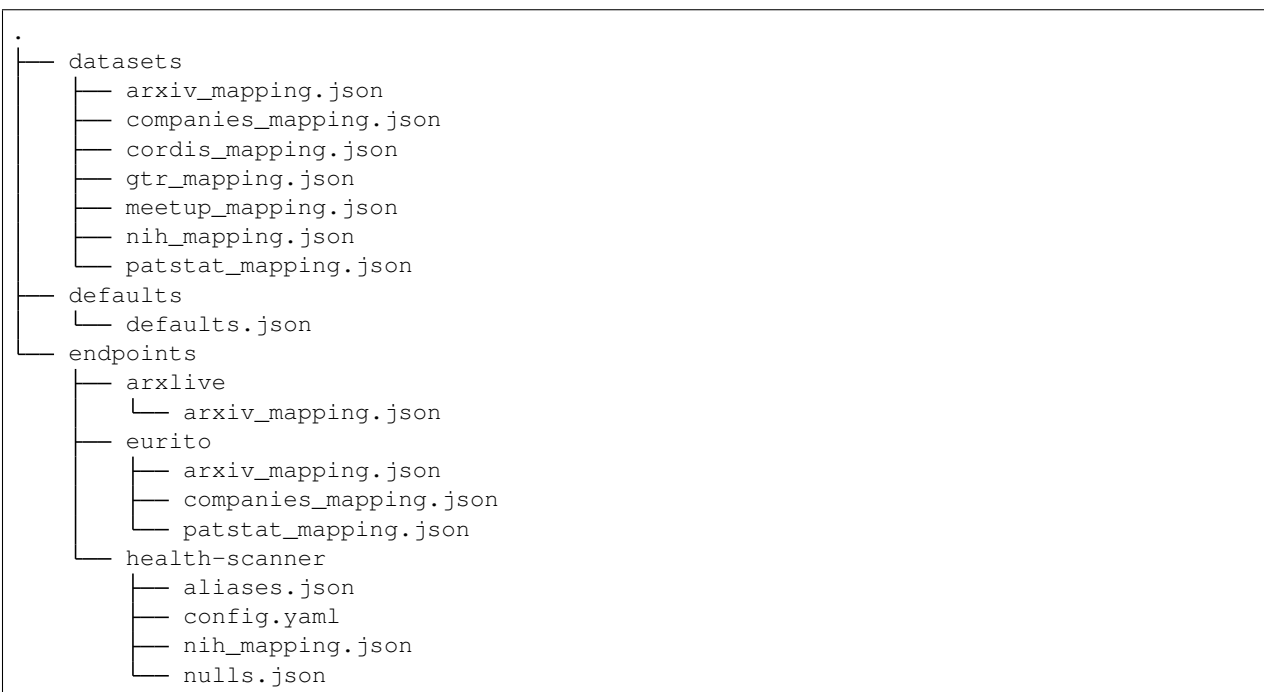
Tier 0 fields are implicitly excluded from tier 1 if they are missing from the `schema_transformation` file. Tier 1 schema field names are applied via `nesta.packages.decorator.schema_transform`

Tier 2

Although not-yet-implemented, the tier 2 schema is reserved for future graph ontologies. Don't expect any changes any time soon!

2.2.5 Elasticsearch mappings

Our methodology for constructing Elasticsearch mappings is described here. It is intended to minimise duplication of efforts and enforce standardisation when referring to a common dataset whilst being flexible to individual project needs. It is implied in our framework that a single dataset can be used across many projects, and each project is mapped to a single endpoint. It is useful to start by looking at the structure of the `nesta/core/schemas/tier_1/mappings/` directory:



Firstly we consider `defaults/defaults.json` which should contain all default fields for all mappings - for example standard analyzers and dynamic strictness. We might also consider putting global fields there.

Next consider the `datasets` subdirectory. Each mapping file in here should contain the complete mappings field for the respective dataset. The naming convention `<dataset>_mapping.json` is a hard requirement, as `<dataset>` will map to the index for this dataset at any given endpoint.

Finally consider the `endpoints` subdirectory. Each sub-subdirectory here should map to any endpoint which requires changes beyond the `defaults` and `datasets` mappings. Each mapping file within each

endpoint sub-subdirectory (e.g. arxlive or health-scanner) should satisfy the same naming convention (<dataset>_mapping.json). All conventions here are also consistent with the `elasticsearch.yaml` configuration file (to see this configuration, you will need to clone the repo and follow [these steps](#) to unencrypt the config), which looks a little like this:

```
## The following assumes the AWS host endpoint naming convention:
## {scheme}://search-{endpoint}-{id}.{region}.es.amazonaws.com
defaults:
  scheme: https
  port: 443
  region: eu-west-2
  type: _doc
endpoints:
  # -----
  # <AWS endpoint domain name>:
  #   id: <AWS endpoint UUID>
  #   <default override key>: <default override value> ## e.g.: scheme, port, region,
  #   _type
  #   indexes:
  #     <index name>: <incremental version number> ## Note: defaults to <index name>_
  # dev in testing mode
  # -----
  arxlive:
    id: <this is a secret>
    indexes:
      arxiv: 4
  # -----
  health-scanner:
    id: <this is a secret>
    indexes:
      nih: 6
      companies: 5
      meetup: 4
  ... etc ...
```

Note that for the health-scanner endpoint, companies and meetup will be generated from the datasets mappings, as they are not specified under the endpoints/health-scanner subdirectory. Also note that endpoints sub-directories do not need to exist for each endpoint to be generated: the mappings will simply be generated from the dataset defaults. For example, a new endpoint general can be generated from the DAPS codebase using the above, even though there is no endpoints/general sub-subdirectory.

Individual endpoints can also specify `aliases.json` which harmonises field names across datasets for specific endpoints. This uses a convention as follows:

```
{
  #...the convention is...
  "<new field name>": {
    "<dataset 1>": "<old field name 1>",
    "<dataset 2>": "<old field name 2>",
    "<dataset 3>": "<old field name 3>"
  },
  #...an example is...
  "city": {
    "companies": "placeName_city_organisation",
    "meetup": "placeName_city_group",
    "nih": "placeName_city_organisation"
  },
  #...etc...#
```

(continues on next page)

(continued from previous page)

```
}
```

By default, this applies (what Joel calls) a “soft” alias, which is an [Elasticsearch alias](#), however by specifying `hard-alias=true` in `config.yaml` (see `health-scanner` above), the alias is instead applied directly (i.e. field names are physically replaced, not aliased).

You will also notice the `nulls.json` file in the `health-scanner` endpoint. This is a relatively experimental feature for automatically nullifying values on ingestion through ElasticsearchPlus, in lieu of proper exploratory data analysis. The logic and format for this is [documented here](#).

Mapping construction hierarchy

Each mapping is constructed by overriding nested fields using the defaults datasets and endpoints, in that order (i.e. endpoints override nested fields in datasets, and datasets override those in defaults). If you would like to “switch off” a field from the defaults or datasets mappings, you should set the value of the nested field to `null`. For example:

```
{
  "mappings": {
    "_doc": {
      "dynamic": "strict",
      "properties": {
        "placeName_zipcode_organisation": null
      }
    }
  }
}
```

will simply “switch off” the field `placeName_zipcode_organisation`, which was specified in datasets.

The logic for the mapping construction hierarchy is demonstrated in the respective `orms.orm_utils.get_es_mapping` function:

```
def get_es_mapping(dataset, endpoint):
    '''Load the ES mapping for this dataset and endpoint,
    including aliases.

    Args:
        dataset (str): Name of the dataset for the ES mapping.
        endpoint (str): Name of the AWS ES endpoint.

    Returns:
        :obj:`dict`
    '''
    mapping = _get_es_mapping(dataset, endpoint)
    _apply_alias(mapping, dataset, endpoint)
    _prune_nested(mapping) # prunes any nested keys with null values
    return mapping
```

Integrated tests

The following `pytest` tests are made (and triggered on PR via travis):

- `aliases.json` files are checked for consistency with available datasets.

- All mappings for each in `datasets` and `endpoints` are fully generated, and tested for compatibility with the schema transformations (which are, in turn, checked against the valid ontology in `ontology.json`).

Features in DAPS2

- The index version (e.g. `'arxiv': 4` in `elasticsearch.yaml`) will be automatically generated from semantic versioning and the git hash in DAPS2, therefore the `indexes` field will consolidate to an itemised list of indexes.
- The mappings under `datasets` will be automatically generated from the open ontology which will be baked into the tier-0 schemas. This will render `schema_transformations` redundant.
- Elasticsearch components will be factored out of `orm_utils`.

2.2.6 Luigi Hacks

Modifications and possible future contributions to the Luigi module.

autobatch

batchclient

NOTE: overwhelmingly based on [this](#), where the following documentation has been directly lifted. The main difference to the latter, is that AWS jobs are submitted via `**kwargs` in order to allow more flexibility (and probably more future-proofing if new parameters are added to `boto3`).

AWS Batch wrapper for Luigi

From the AWS website:

AWS Batch enables you to run batch computing workloads on the AWS Cloud.

Batch computing is a common way for developers, scientists, and engineers to access large amounts of compute resources, and AWS Batch removes the undifferentiated heavy lifting of configuring and managing the required infrastructure. AWS Batch is similar to traditional batch computing software. This service can efficiently provision resources in response to jobs submitted in order to eliminate capacity constraints, reduce compute costs, and deliver results quickly.

See [AWS Batch User Guide](#) for more details.

To use AWS Batch, you create a `jobDefinition` JSON that defines a `docker run` command, and then submit this JSON to the API to queue up the task. Behind the scenes, AWS Batch auto-scales a fleet of EC2 Container Service instances, monitors the load on these instances, and schedules the jobs.

This `boto3-powered` wrapper allows you to create Luigi Tasks to submit Batch `jobDefinition`'s. You can either pass a dict (mapping directly to the `jobDefinition` JSON) OR an Amazon Resource Name (arn) for a previously registered `jobDefinition`.

Requires:

- `boto3` package
- Amazon AWS credentials discoverable by `boto3` (e.g., by using `aws configure` from `awscli`)
- An enabled AWS Batch job queue configured to run on a compute environment.

Written and maintained by Jake Feala (@jfeala) for Outlier Bio (@outlierbio)

exception BatchJobException

Bases: `Exception`

class BatchClient (*poll_time=10, **kwargs*)

Bases: `object`

get_active_queue ()

Get name of first active job queue

get_job_id_from_name (*job_name*)

Retrieve the first job ID matching the given name

get_job_status (*job_id*)

Retrieve task statuses from ECS API

Parameters (*str*) (*job_id*) – AWS Batch job uuid

Returns one of {SUBMITTED|PENDING|RUNNABLE|STARTING|RUNNING|SUCCEEDED|FAILED}

get_logs (*log_stream_name*, *get_last=50*)

Retrieve log stream from CloudWatch

submit_job (***kwargs*)

Wrap submit_job with useful defaults

terminate_job (***kwargs*)

Wrap terminate_job

hard_terminate (*job_ids*, *reason*, *iattempt=0*, ***kwargs*)

Terminate all jobs with a hard(ish) exit via an Exception. The function will also wait for jobs to be explicitly terminated

wait_on_job (*job_id*)

Poll task status until STOPPED

register_job_definition (*json_fpath*)

Register a job definition with AWS Batch, using a JSON

misctools

A collection of miscellaneous tools.

get_config (*file_name*, *header*, *path='core/config'*)

Get the configuration from a file in the luigi config path directory, and convert the key-value pairs under the config header into a *dict*.

Parameters

- **file_name** (*str*) – The configuration file name.
- **header** (*str*) – The header key in the config file.

Returns *dict*

get_paths_from_relative (*relative=1*)

A helper method for within *find_filepath_from_pathstub*. Prints all file and directory paths from a relative number of ‘backward steps’ from the current working directory.

find_filepath_from_pathstub (*path_stub*)

Find the full path of the ‘closest’ file (or directory) to the current working directory ending with *path_stub*. The *closest* file is determined by starting forwards of the current working directory. The algorithm is then repeated by moving the current working directory backwards, one step at a time until the file (or directory) is found. If the HOME directory is reached, the algorithm raises `FileNotFoundError`.

Parameters `path_stub` (*str*) – The partial file (or directory) path stub to find.

Returns The full path to the partial file (or directory) path stub.

f3p (*path_stub*)

Shortened name for coding convenience

load_yaml_from_pathstub (*pathstub*, *filename*)

Basic wrapper around `find_filepath_from_pathstub` which also opens the file (assumed to be yaml).

Parameters

- **pathstub** (*str*) – Stub of filepath where the file should be found.
- **filename** (*str*) – The filename.

Returns The file contents as a json-like object.

load_batch_config (*luigi_task*, *additional_env_files*=[], ***overrides*)

Load default luigi batch parameters, and apply any overrides if required. Note that the usage pattern for this is normally `load_batch_config(self, additional_env_files, **overrides)` from within a luigi Task, where `self` is the luigi Task.

Parameters

- **luigi_task** (*luigi.Task*) – Task to extract test and date parameters from.
- **additional_env_files** (*list*) – List of files to pass directly to the batch local environment.
- **overrides** (***kwargs*) – Any overrides or additional parameters to pass to the batch task as parameters.

Returns Batch configuration parameters, which can be expanded as ***kwargs* in BatchTask.

Return type `config` (*dict*)

extract_task_info

Extract task name and generate a routine id from a luigi task, from the date and test fields.

Parameters **luigi_task** (*luigi.Task*) – Task to extract test and date parameters from.

Returns Test flag, and routine ID for this task.

Return type {*test*, *routine_id*} (*tuple*)

mysqldb

NOTE: overwhelmingly based on [this2](#), where the following documentation has been directly lifted. The main difference to the latter, is that ***cnx_kwargs* in the constructor can accept *port* as a key.

make_mysql_target (*luigi_task*, *mysqldb_env*=*'MYSQLDB'*)

Generate a MySQL target for a luigi Task, based on the Task's date and test parameters, and indicated configuration file.

Parameters

- **luigi_task** (*luigi.Task*) – Task to extract test and date parameters from.
- **mysqldb_env** (*str*) – Environmental variable storing the path to MySQL config.

Returns target (*MySQLTarget*)


```
class MySQLTarget (host, database, user, password, table, update_id, **cnx_kwargs)
    Bases: luigi.target.Target

    Target for a resource in MySQL.

    marker_table = 'table_updates'

    touch (connection=None)
        Mark this update as complete.

        IMPORTANT, If the marker table doesn't exist, the connection transaction will be aborted and the connection reset. Then the marker table will be created.

    exists (connection=None)
        Returns True if the Target exists and False otherwise.

    connect (autocommit=False)

    create_marker_table ()
        Create marker table if it doesn't exist.

        Using a separate connection since the transaction might have to be reset.
```

s3

A more recent implementation of AWS S3 support, stolen from: https://gitlab.com/ced/s3_helpers/blob/master/luigi_s3_target.py, but instead using modern boto3 commands.

```
merge_dicts (*dicts)
    Merge dicts together, with later entries overriding earlier ones.

parse_s3_path (path)
    For a given S3 path, return the bucket and key values

class S3FS (**kwargs)
    Bases: luigi.target.FileSystem

    exists (path)
        Return true if S3 key exists

    remove (path, recursive=True)
        Remove a file or directory from S3

    mkdir (path, parents=True, raise_if_exists=False)
        Create directory at location path

        Creates the directory at path and implicitly create parent directories if they do not already exist.

        Parameters

        • path (str) – a path within the FileSystem to create as a directory.

        • parents (bool) – Create parent directories when necessary. When parents=False and the parent directory doesn't exist, raise luigi.target.MissingParentDirectory

        • raise_if_exists (bool) – raise luigi.target.FileAlreadyExists if the folder already exists.

    isdir (path)
        Return True if the location at path is a directory. If not, return False.

        Parameters path (str) – a path within the FileSystem to check as a directory.

    Note: This method is optional, not all FileSystem subclasses implements it.
```

listdir (*path*)

Return a list of files rooted in *path*.

This returns an iterable of the files rooted at *path*. This is intended to be a recursive listing.

Parameters *path* (*str*) – a path within the FileSystem to list.

Note: This method is optional, not all FileSystem subclasses implements it.

copy (*path*, *dest*)

Copy a file or a directory with contents. Currently, LocalFileSystem and MockFileSystem support only single file copying but S3Client copies either a file or a directory as required.

move (*path*, *dest*)

Move a file, as one would expect.

du (*path*)

class S3Target (*path*, *s3_args*={}, ***kwargs*)

Bases: `luigi.target.FileSystemTarget`

fs = None

open (*mode*='rb')

Open the FileSystem target.

This method returns a file-like object which can either be read from or written to depending on the specified mode.

Parameters *mode* (*str*) – the mode *r* opens the FileSystemTarget in read-only mode, whereas *w* will open the FileSystemTarget in write mode. Subclasses can implement additional options.

class AtomicS3File (*path*, *s3_obj*, ***kwargs*)

Bases: `luigi.target.AtomicLocalFile`

move_to_final_destination ()

2.2.7 Scripts

A set of helper scripts for the batching system.

Note that this directory is required to sit in *\$PATH*. By convention, all executables in this directory start with *nesta_* so that our developers know where to find them.

nesta_prepare_batch

Collect a batchable `run.py` file, including dependencies and an automatically generated requirements file; which is all zipped up and sent to AWS S3 for batching. This script is executed automatically in `luigi hacks.autobatch.AutoBatchTask.run`.

Parameters:

- **BATCHABLE_DIRECTORY:** The path to the directory containing the batchable `run.py` file.
- **ARGS:** Space-separated-list of files or directories to include in the zip file, for example imports.

nesta_docker_build

Build a docker environment and register it with the AWS ECS container repository.

Parameters:

- **DOCKER_RECIPE:** A docker recipe. See `docker_recipes/` for a good idea of how to build a new environment.

2.2.8 Elasticsearch

The following steps will take you through setting up elasticsearch on an EC2 instance.

Launch the EC2 instance and ssh in so the following can be installed:

docker

```
sudo yum install docker -y
```

docker-compose

```
curl -L https://github.com/docker/compose/releases/download/1.22.0/
docker-compose-`uname -s` - `uname -m` -o /usr/local/bin/docker-compose
chmod +x /usr/local/bin/docker-compose          sudo ln -s /usr/local/bin/
docker-compose /usr/bin/docker-compose
```

more info: <https://github.com/docker/compose/releases>

docker permissions

```
sudo usermod -a -G docker $USER
```

more info: <https://techoverflow.net/2017/03/01/solving-docker-permission-denied-while-trying-to-connect-to-the-docker-daemon-sock>

vm.max_map_count

set permanantly in `/etc/sysctl.conf` by adding the following line: `vm.max_map_count=262144`

more info: <https://www.elastic.co/guide/en/elasticsearch/reference/current/docker.html>

python 3.6

```
sudo yum install python36 -y
```

The machine now needs to be rebooted `sudo reboot now`

Docker

the `docker-compose.yml` needs to include ulimits settings::

ulimits:

```
memlock: soft: -1 hard: -1
```

nofile: soft: 65536 hard: 65536

Recipes for http or https clusters can be found in: `nesta/core/scripts/elasticsearch`

There is also an EC2 AMI for a http node stored in the London region: `elasticsearch node London vX`

Reindexing data from a remote cluster

- reindex permissions need to be set in the new cluster's `elasticsearch.yml`
- if the existing cluster is AWS hosted ES the ip address needs to be added to the security settings
- follow this guide: <https://www.elastic.co/guide/en/elasticsearch/reference/current/docs-reindex.html#reindex-from-remote>
- `index` and `query` do not need to be supplied
- if reindexing from AWS ES the port should be 443 for https. This is mandatory in the json sent to the reindexing api

2.2.9 Containerised Luigi

Requirements

To containerise a pipeline a few steps are required:

- All imports must be absolute, ie `from nesta . packages, core` etc
- Once testing is complete the code should be committed and pushed to github, as this prevents the need to use local build options
- If building and running locally, Docker must be installed on the local machine and given enough RAM in the settings to run the pipeline
- Any required configuration files must be in `nesta.core.config` ie luigi and mysql config files, any API keys. **This directory is ignored but check before committing**

Build

The build uses a multi-stage Dockerfile to reduce the size of the final image:

1. Code is git cloned and requirements are pip installed into a virtual environment
2. The environment is copied into the second image

From the root of the repository: `docker build -f docker/Dockerfile -t name:tag .`

Where `name` is the name of the created image and `tag` is the chosen tag. Eg `arxlive:latest`. This just makes the run step easier rather than using the generated id

The two stage build will normally just rebuild the second stage pulling in new code only. If a full rebuild is required, eg after `requirements.txt` has changed then include: `--no-cache`

Python version defaults to 3.6 but can be set during build by including the flag: `--build-arg python-version=3.7`

Tag defaults to `dev` but this can be overridden by including the flag: `--build-arg GIT_TAG=0.3` a branch name also works `--build-arg GIT_TAG=my-branch`

Work from a branch or locally while testing. Override the target branch from Github using the method above, or use the commented out code in the Dockerfile to switch to build from local files. **Don't commit this change!**

Run

As only one pipeline runs in the container the `luigid` scheduler is not used.

There is a `docker-compose` file for arXlive which mounts your local `~.aws` folder for AWS credentials as this outside docker's context:

```
docker-compose -f docker/docker-compose.yml run luigi --module module_path
params
```

Where:

- `docker-compose.yml` is the docker-compose file containing the image: `image_name:tag` from the build
- `module_path` is the full python path to the module
- `params` are any other params to supply as per normal, ie `--date --production` etc

Eg: `docker-compose -f docker/docker-compose-arxlive-dev.yml run luigi --module nesta.core.routines.arxiv.arxiv_iterative_root_task RootTask --date 2019-04-16`

This could be adapted for each pipeline, or alternatively run with the volume specified with `-v`

```
docker run -v ~/.aws:/root/.aws:ro name:tag --module ...
```

Where `name` is the name of the created image and `tag` is the chosen tag. Eg `arxlive:latest` `--module ...` onwards contains the arguments you would pass to Luigi.

Scheduling

1. Create an executable shell script in `nesta.core.scripts` to launch docker-compose with all the necessary parameters eg: `production`
2. Add a cron job to the shell script (there are some good online cron syntax checking tools, if needed)
3. Set the cron job to run every few minutes while testing and check the logs with `docker logs mycontainerterhash --tail 50`. Obtain the hash using `docker ps`
4. It will run logged in as the user who set it up but there still may still be some permissions issues to deal with

Currently scheduled

arXlive:

- A shell script to launch docker-compose for arXlive is set up to run in a cron job on user `russellwinch`
- This is scheduled for Sunday-Thursday at 0300 GMT. arXiv is updated on these days at 0200 GMT
- Logs are just stored in the container, use `docker logs container_name` to view

Important points

- Keep any built images secure, they contain credentials
- You need to rebuild if code has changed
- As there is no central scheduler there is nothing stopping you from running the task more than once at the same time, by launching the container multiple times
- The graphical interface is not enabled without the scheduler

Debugging

If necessary, it's possible to debug inside the container, but the endpoint needs to be overridden with bash:

```
docker run --entrypoint /bin/bash -itv ~/.aws:/root/.aws:ro image_name:tag
```

Where `image_name:tag` is the image from the build step This includes the mounting of the `.aws` folder

Almost nothing is installed (not even vi!!) other than Python so `apt-get update` and then `apt-get install` whatever you need

Todo

A few things are sub-optimal:

- The container runs on the prod box rather than in the cloud in ECS
- Credentials are held in the container and local AWS config is required, this is the cause of the above point
- Due to the Nesta monorepo everything is pip installed, making a large container size with many unrequired packages. Pipeline specific requirements should be considered
- As logs are stored in the old containers they are kept until the next run where they are pruned and the logs are lost. Add a method of getting the logs to the host logger and record centrally
- In the arXlive pipeline there are at least 500 calls to the MAG API each run as the process tries to pick up new title matches on existing articles. As the API key only allows 10,000 calls per month this is currently OK with the schedule as it is but could possibly go over at some point

3.1 Where is the data?

As a general rule-of-thumb, our data is always stored in the London region (`eu-west-2`), in either RDS (tier-0, MySQL) or Elasticsearch (tier-1). For the EURITO project we also use Neo4j (tier-1), and in the distant future we will use Neo4j for tier-2 (i.e. a knowledge graph).

3.2 Why don't you use Aurora rather than MySQL?

Aurora is definitely cheaper for stable production and business processes but not for research and development. You are charged for every byte of data you have *ever* consumed. This quickly spirals out-of-control for big-data development. Maybe one day we'll consider migrating back, once the situation stabilises.

3.3 Where are the production machines?

Production machines (EC2) run in Ohio (`us-east-2`).

3.4 Where is the latest config?

We use `git-crypt` to encrypt our configuration files whilst allowing them to be versioned in git (meaning that we can also rollback configuration). To unlock the configuration encryption, you should install `git-crypt`, then run `bash install.sh` from the project root, and finally unlock the configuration using the key found [here](#).

3.5 Where do I start with Elasticsearch?

All Elasticsearch indexes (aka “databases” to the rest of the world), mappings (aka “schemas”) and whitelisting can be found [here](#).

I’d recommend using PostMan for spinning up and knocking down indexes. Practice this on a new cluster (which you can spin up from the above link), and then practice `PUT`, `POST` and `DELETE` requests to `PUT` an index (remember: “database”) with a mapping (“schema”), inserting a “row” with `POST` and then deleting the index with `DELETE`. You will quickly learn that it’s very easy to delete everything in Elasticsearch.

4.1 I'm having problems using the config files!

We use `git-crypt` to encrypt our configuration files whilst allowing them to be versioned in git (meaning that we can also rollback configuration). To unlock the configuration encryption, you should install `git-crypt`, then run `bash install.sh` from the project root, and finally unlock the configuration using the [key](#).

4.2 How do I restart the apache server after downtime?

```
sudo service httpd restart
```

4.3 How do I restart the luigi server after downtime?

```
sudo su - luigi
source activate py36
luigid --background --pidfile /var/run/luigi/luigi.pid --logdir /var/log/luigi
```

4.4 How do I perform initial setup to ensure the batchables will run?

- AWS CLI needs to be installed and configured:

```
pip install awscli
aws configure
```

AWS Access Key ID and Secret Access Key are set up in IAM > Users > Security Credentials Default region name should be `eu-west-1` to enable the error emails to be sent In AWS SES the sender and receiver email addresses need to be verified

- The config files need to be accessible and the PATH and LUIGI_CONFIG_PATH need to be amended accordingly

4.5 How can I send/receive emails from Luigi?

You should set the environmental variable `export LUIGI_EMAIL="<your.email@something>"` in your `.bashrc`. You can test this with `luigi TestNotificationsTask --local-scheduler --email-force-send`. Make sure your email address has been registered under AWS SES.

4.6 How do I add a new user to the server?

- add the user with `useradd --create-home username`
- add sudo privileges [following these instructions](#)
- add to ec2 user group with `sudo usermod -a -G ec2-user username`
- set a temp password with `passwd username`
- their home directory will be `/home/username/`
- copy `.bashrc` to their home directory
- create folder `.ssh` in their home directory
- copy `.ssh/authorized_keys` to the same folder in their home directory (DONT MOVE IT!!)
- `cd` to their home directory and perform the below
- `chown` their copy of `.ssh/authorized_keys` to their username: `chown username .ssh/authorized_keys`
- clone the nesta repo
- copy `core/config` files
- set password to be changed next login `chage -d 0 username`
- share the temp password and core pem file

If necessary: `- sudo chmod g+w /var/tmp/batch`

CHAPTER 5

Packages

Nesta's collection of tools for meaty tasks. Any processes that go into production come here first, but there are other good reasons for code to end up here.

Nesta's production system is based on [Luigi](#) pipelines, and are designed to be entirely run on AWS via the batch service. The main Luigi server runs on a persistent EC2 instance. Beyond the well documented Luigi code, the main features of the nesta production system are:

- `luigihacks.autobatch`, which facilitates a managed `Luigi.Task` which is split, batched and combined in a single step. Currently only synchronous jobs are accepted. Asynchronous jobs (where downstream `Luigi.Task` jobs can be triggered) are a part of a longer term plan.
- `scripts.nesta_prepare_batch` which zips up the batchable with the specified environmental files and ships it to AWS S3.
- `scripts.nesta_docker_build` which builds a specified docker environment and ships it to AWS ECS.

6.1 License

MIT © 2018 Nesta

Python Module Index

c

core, 13
core.batchables, 22
core.batchables.batchgeocode.run, 25
core.batchables.cordis.cordis_api.run, 22
core.batchables.crunchbase.crunchbase_collect.run, 23
core.batchables.examples.batch_example.run, 22
core.batchables.examples.template_batchable.run, 22
core.batchables.meetup.country_groups.run, 25
core.batchables.meetup.group_details.run, 25
core.batchables.meetup.groups_members.run, 25
core.batchables.meetup.members_groups.run, 25
core.batchables.meetup.topic_tag_elasticsearch.run, 25
core.batchables.nih.nih_abstract_mesh_data.run, 24
core.batchables.nih.nih_dedupe.run, 24
core.batchables.nih.nih_process_data.run, 24
core.luigihacks, 34
core.luigihacks.batchclient, 34
core.luigihacks.mistools, 35
core.luigihacks.mysqlldb, 36
core.luigihacks.s3, 37
core.orms.meetup_orm, 27
core.orms.nih_orm, 28
core.routines, 14
core.routines.examples.db_example.db_example.run, 15
core.routines.examples.s3_example.s3_example.run, 14

n

nesta.core.batchables.eurito.patstat_eu.run, 23
nesta.core.batchables.gtr.collect_gtr.run, 23
nesta.core.batchables.gtr.embed_topics.run, 24
nesta.core.batchables.nlp.corex_topic_model.run, 26
nesta.core.batchables.nlp.ngrammer.run, 26
nesta.core.batchables.nlp.tfidf.run, 26
nesta.core.batchables.nlp.vectorizer.run, 26
nesta.core.batchables.novelty.lolvelty.run, 27
nesta.core.routines.arxiv.arxiv_collect_iterative_task.run, 16
nesta.core.routines.arxiv.arxiv_grid_task, 19
nesta.core.routines.arxiv.arxiv_iterative_date_task.run, 16
nesta.core.routines.arxiv.arxiv_mag_sparql_task, 18
nesta.core.routines.arxiv.arxiv_mag_task, 17
nesta.core.routines.meetup.health_tagging.topic_discovery.run, 21

p

packages, 3
packages.decorators.ratelimit, 10
packages.decorators.schema_transform, 10
packages.format_utils.datetools, 9
packages.geo_utils.geocode, 8
packages.meetup.country_groups, 3
packages.meetup.group_details, 6
packages.meetup.groups_members, 5

`packages.meetup.meetup_utils`, 6
`packages.meetup.members_groups`, 5
`packages.nih.collect_nih`, 7
`packages.nih.process_nih`, 7
`packages.nlp_utils.preprocess`, 8

A

abstract_text (*Abstracts attribute*), 29
 Abstracts (*class in core.orms.nih_orm*), 29
 activity (*Projects attribute*), 28
 administering_ic (*Projects attribute*), 28
 affiliation (*Publications attribute*), 29
 application_id (*Abstracts attribute*), 29
 application_id (*Projects attribute*), 28
 application_type (*Projects attribute*), 28
 arra_funded (*Projects attribute*), 28
 articles_from_date (*CollectNewTask attribute*), 16
 articles_from_date (*DateTask attribute*), 17
 articles_from_date (*GridRootTask attribute*), 20
 articles_from_date (*GridTask attribute*), 20
 articles_from_date (*MagSparqlTask attribute*), 19
 articles_from_date (*QueryMagTask attribute*), 18
 assert_iso2_key() (*in module packages.meetup.country_groups*), 4
 AtomicS3File (*class in core.luigihacks.s3*), 38
 author_list (*Publications attribute*), 29
 author_name (*Publications attribute*), 29
 award_notice_date (*Projects attribute*), 28

B

BatchClient (*class in core.luigihacks.batchclient*), 35
 BatchJobException, 34
 budget_end (*Projects attribute*), 28
 budget_start (*Projects attribute*), 28

C

category_id (*Group attribute*), 27
 category_name (*Group attribute*), 27
 category_shortcode (*Group attribute*), 27
 cfda_code (*Projects attribute*), 28
 chunker() (*in module nesta.core.batchables.nlp.tfidf.run*), 26
 city (*Group attribute*), 27

clean_abstract() (*in module core.batchables.nih.nih_abstract_mesh_data.run*), 24
 clean_and_tokenize() (*in module packages.nlp_utils.preprocess*), 8
 ClinicalStudies (*class in core.orms.nih_orm*), 30
 clinicaltrials_gov_id (*ClinicalStudies attribute*), 30
 CollectNewTask (*class in nesta.core.routines.arxiv.arxiv_collect_iterative_task*), 16
 connect() (*MySQLTarget method*), 37
 copy() (*S3FS method*), 38
 core (*module*), 13
 core.batchables (*module*), 22
 core.batchables.batchgeocode.run (*module*), 25
 core.batchables.cordis.cordis_api.run (*module*), 22
 core.batchables.crunchbase.crunchbase_collect.run (*module*), 23
 core.batchables.examples.batch_example.run (*module*), 22
 core.batchables.examples.template_batchable.run (*module*), 22
 core.batchables.meetup.country_groups.run (*module*), 25
 core.batchables.meetup.group_details.run (*module*), 25
 core.batchables.meetup.groups_members.run (*module*), 25
 core.batchables.meetup.members_groups.run (*module*), 25
 core.batchables.meetup.topic_tag_elasticsearch.run (*module*), 25
 core.batchables.nih.nih_abstract_mesh_data.run (*module*), 24
 core.batchables.nih.nih_dedupe.run (*module*), 24
 core.batchables.nih.nih_process_data.run

(*module*), 24
 core.luigihacks (*module*), 34
 core.luigihacks.batchclient (*module*), 34
 core.luigihacks.misctools (*module*), 35
 core.luigihacks.mysqlldb (*module*), 36
 core.luigihacks.s3 (*module*), 37
 core.orms.meetup_orm (*module*), 27
 core.orms.nih_orm (*module*), 28
 core.routines (*module*), 14
 core.routines.examples.db_example.db_example (*module*), 15
 core.routines.examples.s3_example.s3_example (*module*), 14
 core_categories (*TopicDiscoveryTask* attribute), 21
 core_project_num (*Projects* attribute), 28
 core_project_number (*ClinicalStudies* attribute), 30
 country (*Group* attribute), 27
 country (*Publications* attribute), 29
 country_code (*MeetupCountryGroups* attribute), 4
 country_name (*Group* attribute), 27
 create_marker_table() (*MySQLTarget* method), 37
 created (*Group* attribute), 27

D

date (*CollectNewTask* attribute), 16
 date (*DateTask* attribute), 17
 date (*FinalTask* attribute), 14
 date (*GridRootTask* attribute), 20
 date (*GridTask* attribute), 19
 date (*InputData* attribute), 15
 date (*MagSparqlTask* attribute), 18
 date (*QueryMagTask* attribute), 18
 date (*RootTask* attribute), 16
 date (*SomeTask* attribute), 14, 15
 DateTask (*class* in *nesta.core.routines.arxiv.arxiv_iterative_date_task*), 17
 db_config (*InputData* attribute), 15
 db_config_env (*CollectNewTask* attribute), 16
 db_config_env (*DateTask* attribute), 17
 db_config_env (*GridTask* attribute), 19
 db_config_env (*MagSparqlTask* attribute), 19
 db_config_env (*QueryMagTask* attribute), 18
 db_config_env (*TopicDiscoveryTask* attribute), 21
 db_config_path (*CollectNewTask* attribute), 16
 db_config_path (*DateTask* attribute), 17
 db_config_path (*GridRootTask* attribute), 20
 db_config_path (*GridTask* attribute), 19
 db_config_path (*MagSparqlTask* attribute), 19
 db_config_path (*QueryMagTask* attribute), 18
 debug (*GridRootTask* attribute), 20
 description (*Group* attribute), 27
 direct_cost_amt (*Projects* attribute), 29

drop_and_recreate (*GridRootTask* attribute), 20
 du() (*S3FS* method), 38

E

ed_inst_type (*Projects* attribute), 28
 exists() (*MySQLTarget* method), 37
 exists() (*S3FS* method), 37
 extract_core_orgs() (in *module* *core.batchables.cordis.cordis_api.run*), 22
 extract_date() (in *module* *packages.format_utils.datetools*), 10
 extract_task_info (in *module* *core.luigihacks.misctools*), 36
 extract_year() (in *module* *packages.format_utils.datetools*), 10
 extract_yearly_funds() (in *module* *core.batchables.nih.nih_dedupe.run*), 24

F

f3p() (in *module* *core.luigihacks.misctools*), 36
 filter_by_idf() (in *module* *packages.nlp_utils.preprocess*), 8
 FinalTask (*class* in *core.routines.examples.s3_example.s3_example*), 14
 find_filepath_from_pathstub() (in *module* *core.luigihacks.misctools*), 35
 flatten_data() (in *module* *packages.meetup.meetup_utils*), 6
 foa_number (*Projects* attribute), 28
 fs (*S3Target* attribute), 38
 full_project_num (*Projects* attribute), 28
 funding_ics (*Projects* attribute), 28
 funding_mechanism (*Projects* attribute), 28
 fy (*Projects* attribute), 28

G

generate_composite_key() (in *module* *packages.geo_utils.geocode*), 9
 generate_coords() (in *module* *packages.meetup.country_groups*), 3
 geocode() (in *module* *packages.geo_utils.geocode*), 8
 geocode_batch_dataframe() (in *module* *packages.geo_utils.geocode*), 9
 geocode_dataframe() (in *module* *packages.geo_utils.geocode*), 9
 get_active_queue() (*BatchClient* method), 35
 get_all_members() (in *module* *packages.meetup.groups_members*), 5
 get_api_key() (in *module* *packages.meetup.meetup_utils*), 6
 get_config() (in *module* *core.luigihacks.misctools*), 35

[get_coordinate_data\(\)](#) (in module *pack-ages.meetup.country_groups*), 4
[get_core_topics\(\)](#) (in module *pack-ages.meetup.meetup_utils*), 7
[get_data_urls\(\)](#) (in module *pack-ages.nih.collect_nih*), 7
[get_group_details\(\)](#) (in module *pack-ages.meetup.group_details*), 6
[get_groups\(\)](#) (*MeetupCountryGroups* method), 4
[get_groups_recursive\(\)](#) (*MeetupCountryGroups* method), 5
[get_job_id_from_name\(\)](#) (*BatchClient* method), 35
[get_job_status\(\)](#) (*BatchClient* method), 35
[get_logs\(\)](#) (*BatchClient* method), 35
[get_member_details\(\)](#) (in module *pack-ages.meetup.members_groups*), 5
[get_member_groups\(\)](#) (in module *pack-ages.meetup.members_groups*), 5
[get_members\(\)](#) (in module *pack-ages.meetup.groups_members*), 5
[get_members_by_percentile\(\)](#) (in module *packages.meetup.meetup_utils*), 6
[get_paths_from_relative\(\)](#) (in module *core.luigihacks.miscutils*), 35
[get_value\(\)](#) (in module *core.batchables.nih.nih_dedupe.run*), 24
[GridRootTask](#) (class in *nesta.core.routines.arxiv.arxiv_grid_task*), 20
[GridTask](#) (class in *nesta.core.routines.arxiv.arxiv_grid_task*), 19
[Group](#) (class in *core.orms.meetup_or*), 27
[group_id](#) (*GroupMember* attribute), 28
[group_urlname](#) (*GroupMember* attribute), 28
[GroupMember](#) (class in *core.orms.meetup_or*), 28
[groups](#) (*MeetupCountryGroups* attribute), 4

H

[hard_terminate\(\)](#) (*BatchClient* method), 35

I

[ic_name](#) (*Projects* attribute), 28
[id](#) (*Group* attribute), 27
[in_db_config](#) (*SomeTask* attribute), 15
[indirect_cost_amt](#) (*Projects* attribute), 29
[InputData](#) (class in *core.routines.examples.db_example.db_example*), 15
[InputData](#) (class in *core.routines.examples.s3_example.s3_example*), 14
[insert_batch_size](#) (*CollectNewTask* attribute), 16
[insert_batch_size](#) (*DateTask* attribute), 17
[insert_batch_size](#) (*GridRootTask* attribute), 20
[insert_batch_size](#) (*GridTask* attribute), 19
[insert_batch_size](#) (*MagSparqlTask* attribute), 19
[insert_batch_size](#) (*QueryMagTask* attribute), 18
[isdir\(\)](#) (*S3FS* method), 37
[issn](#) (*Publications* attribute), 29
[iterrows\(\)](#) (in module *packages.nih.collect_nih*), 7

J

[journal_issue](#) (*Publications* attribute), 29
[journal_title](#) (*Publications* attribute), 29
[journal_title_abbr](#) (*Publications* attribute), 29
[journal_volume](#) (*Publications* attribute), 29

L

[lang](#) (*Publications* attribute), 30
[lat](#) (*Group* attribute), 27
[LinkTables](#) (class in *core.orms.nih_or*), 30
[listdir\(\)](#) (*S3FS* method), 37
[load_batch_config\(\)](#) (in module *core.luigihacks.miscutils*), 36
[load_transformer\(\)](#) (in module *pack-ages.decorators.schema_transform*), 10
[load_yaml_from_pathstub\(\)](#) (in module *core.luigihacks.miscutils*), 36
[lon](#) (*Group* attribute), 27

M

[mag_config_path](#) (*GridTask* attribute), 19
[mag_config_path](#) (*MagSparqlTask* attribute), 19
[mag_config_path](#) (*QueryMagTask* attribute), 18
[MagSparqlTask](#) (class in *nesta.core.routines.arxiv.arxiv_mag_sparql_task*), 18
[make_mysql_target\(\)](#) (in module *core.luigihacks.mysql*), 36
[marker_table](#) (*MySQLTarget* attribute), 37
[max_age](#) (*SomeTask* attribute), 15
[MeetupCountryGroups](#) (class in *pack-ages.meetup.country_groups*), 4
[member_id](#) (*GroupMember* attribute), 28
[members](#) (*Group* attribute), 27
[members_perc](#) (*TopicDiscoveryTask* attribute), 21
[merge_dicts\(\)](#) (in module *core.luigihacks.s3*), 37
[merge_lists\(\)](#) (in module *nesta.core.batchables.nlp.vectorizer.run*), 27
[metadata\(\)](#) (in module *nesta.core.batchables.eurito.patstat_eu.run*), 23
[mkdir\(\)](#) (*S3FS* method), 37
[move\(\)](#) (*S3FS* method), 38
[move_to_final_destination\(\)](#) (*AtomicS3File* method), 38

MySQLTarget (*class in core.luigihacks.mysqlldb*), 36

N

name (*Group attribute*), 27

nesta.core.batchables.eurito.patstat_eu.run
(*module*), 23

nesta.core.batchables.gtr.collect_gtr.run
(*module*), 23

nesta.core.batchables.gtr.embed_topics.run
(*module*), 24

nesta.core.batchables.nlp.corex_topic_modeling.run
(*module*), 26

nesta.core.batchables.nlp.ngrammer.run
(*module*), 26

nesta.core.batchables.nlp.tfidf.run
(*module*), 26

nesta.core.batchables.nlp.vectorizer.run
(*module*), 26

nesta.core.batchables.novelty.lolvelty.run
(*module*), 27

nesta.core.routines.arxiv.arxiv_collect_patents.run
(*module*), 16

nesta.core.routines.arxiv.arxiv_grid_task.run
(*module*), 19

nesta.core.routines.arxiv.arxiv_iterative_ranking.run
(*module*), 16

nesta.core.routines.arxiv.arxiv_mag_sparql_task.run
(*module*), 18

nesta.core.routines.arxiv.arxiv_mag_task.run
(*module*), 17

nesta.core.routines.meetup.health_tagging.run
(*module*), 21

nih_spending_cats (*Projects attribute*), 29

NoGroupFound, 6

NoMemberFound, 5

O

open () (*S3Target method*), 38

optional () (*in module
nesta.core.batchables.nlp.vectorizer.run*),
27

org_city (*Projects attribute*), 28

org_country (*Projects attribute*), 28

org_dept (*Projects attribute*), 28

org_district (*Projects attribute*), 28

org_duns (*Projects attribute*), 28

org_fips (*Projects attribute*), 28

org_ipf_code (*Projects attribute*), 28

org_name (*Projects attribute*), 28

org_state (*Projects attribute*), 29

org_zipcode (*Projects attribute*), 29

out_db_config (*SomeTask attribute*), 15

output () (*CollectNewTask method*), 16

output () (*FinalTask method*), 14

output () (*GridTask method*), 20

output () (*InputData method*), 14, 15

output () (*MagSparqlTask method*), 19

output () (*QueryMagTask method*), 18

output () (*SomeTask method*), 14, 15

output () (*TopicDiscoveryTask method*), 21

P

packages (*module*), 3

packages.decorators.ratelimit (*module*), 10

packages.decorators.schema_transform
(*module*), 10

packages.format_utils.datetools (*module*),
9

packages.geo_utils.geocode (*module*), 8

packages.meetup.country_groups (*module*), 3

packages.meetup.group_details (*module*), 6

packages.meetup.groups_members (*module*), 5

packages.meetup.meetup_utils (*module*), 6

packages.meetup.members_groups (*module*), 5

packages.events.collect_nih (*module*), 7

packages.nih.process_nih (*module*), 7

packages.nlp_utils.preprocess (*module*), 8

page_number (*Publications attribute*), 30

path () (*in module core.luigihacks.s3*), 37

patent_id (*Patents attribute*), 30

patent_korg_name (*Patents attribute*), 30

patent_title (*Patents attribute*), 30

Patents (*class in core.orms.nih_or*), 30

phr (*Projects attribute*), 29

pi_names (*Projects attribute*), 29

pmc_id (*Publications attribute*), 30

pmid (*LinkTables attribute*), 30

pmid (*Publications attribute*), 29

prepare_data () (*in module
core.batchables.cordis.cordis_api.run*), 22

production (*GridRootTask attribute*), 20

program_officer_name (*Projects attribute*), 29

project_end (*Projects attribute*), 29

project_id (*Patents attribute*), 30

project_number (*LinkTables attribute*), 30

project_start (*Projects attribute*), 29

project_terms (*Projects attribute*), 29

project_title (*Projects attribute*), 29

Projects (*class in core.orms.nih_or*), 28

pub_date (*Publications attribute*), 30

pub_title (*Publications attribute*), 30

pub_year (*Publications attribute*), 30

Publications (*class in core.orms.nih_or*), 29

Q

QueryMagTask (*class in
nesta.core.routines.arxiv.arxiv_mag_task*),

17

R

[ratelimit\(\)](#) (in module [pack-ages.decorators.ratelimit](#)), 10
[register_job_definition\(\)](#) ([BatchClient](#) method), 35
[remove\(\)](#) ([S3FS](#) method), 37
[requires\(\)](#) ([DateTask](#) method), 17
[requires\(\)](#) ([FinalTask](#) method), 14
[requires\(\)](#) ([GridRootTask](#) method), 20
[requires\(\)](#) ([GridTask](#) method), 20
[requires\(\)](#) ([MagSparqlTask](#) method), 19
[requires\(\)](#) ([QueryMagTask](#) method), 18
[requires\(\)](#) ([RootTask](#) method), 16
[requires\(\)](#) ([SomeTask](#) method), 14, 15
[retry_if_not_value_error\(\)](#) (in module [pack-ages.geo_utils.geocode](#)), 9
[RootTask](#) (class in [core.routines.examples.db_example.db_example](#)), 15
[routine_id](#) ([TopicDiscoveryTask](#) attribute), 21
[run\(\)](#) ([CollectNewTask](#) method), 16
[run\(\)](#) ([FinalTask](#) method), 15
[run\(\)](#) ([GridTask](#) method), 20
[run\(\)](#) (in module [core.batchables.batchgeocode.run](#)), 26
[run\(\)](#) (in module [core.batchables.cordis.cordis_api.run](#)), 23
[run\(\)](#) (in module [core.batchables.crunchbase.crunchbase_collect.run](#)), 23
[run\(\)](#) (in module [core.batchables.examples.batch_example.run](#)), 22
[run\(\)](#) (in module [core.batchables.examples.template_batchable.run](#)), 22
[run\(\)](#) (in module [core.batchables.meetup.country_groups.run](#)), 25
[run\(\)](#) (in module [core.batchables.meetup.group_details.run](#)), 25
[run\(\)](#) (in module [core.batchables.meetup.groups_members.run](#)), 25
[run\(\)](#) (in module [core.batchables.meetup.members_groups.run](#)), 25
[run\(\)](#) (in module [core.batchables.meetup.topic_tag_elasticsearch.run](#)), 25
[run\(\)](#) (in module [core.batchables.nih.nih_abstract_mesh_data.run](#)), 24
[run\(\)](#) (in module [core.batchables.nih.nih_dedupe.run](#)), 25
[run\(\)](#) (in module [core.batchables.nih.nih_process_data.run](#)), 24
[run\(\)](#) (in module [nesta.core.batchables.eurito.patstat_eu.run](#)), 23
[run\(\)](#) (in module [nesta.core.batchables.gtr.collect_gtr.run](#)), 24

[run\(\)](#) (in module [nesta.core.batchables.gtr.embed_topics.run](#)), 24
[run\(\)](#) (in module [nesta.core.batchables.nlp.corex_topic_model.run](#)), 26
[run\(\)](#) (in module [nesta.core.batchables.nlp.ngrammer.run](#)), 26
[run\(\)](#) (in module [nesta.core.batchables.nlp.tfidf.run](#)), 26
[run\(\)](#) (in module [nesta.core.batchables.nlp.vectorizer.run](#)), 27
[run\(\)](#) (in module [nesta.core.batchables.novelty.lolvelty.run](#)), 27
[run\(\)](#) ([InputData](#) method), 15
[run\(\)](#) ([MagSparqlTask](#) method), 19
[run\(\)](#) ([QueryMagTask](#) method), 18
[run\(\)](#) ([SomeTask](#) method), 14, 15
[run\(\)](#) ([TopicDiscoveryTask](#) method), 21
[S3FS](#) (class in [core.luigihacks.s3](#)), 37
[S3Target](#) (class in [core.luigihacks.s3](#)), 38
[save_sample\(\)](#) (in module [pack-ages.meetup.meetup_utils](#)), 6
[schema_transform\(\)](#) (in module [pack-ages.decorators.schema_transform](#)), 10
[schema_transformer\(\)](#) (in module [pack-ages.decorators.schema_transform](#)), 10
[select_text\(\)](#) (in module [nesta.core.batchables.eurito.patstat_eu.run](#)), 23
[serial_number](#) ([Projects](#) attribute), 29
[SomeTask](#) (class in [core.routines.examples.db_example.db_example](#)), 15
[SomeTask](#) (class in [core.routines.examples.s3_example.s3_example](#)), 14
[split_links\(\)](#) (in module [core.batchables.cordis.cordis_api.run](#)), 23
[study](#) ([ClinicalStudies](#) attribute), 30
[study_section](#) ([Projects](#) attribute), 29
[study_section_name](#) ([Projects](#) attribute), 29
[study_status](#) ([ClinicalStudies](#) attribute), 30
[submit_job\(\)](#) ([BatchClient](#) method), 35
[subject_id](#) ([Projects](#) attribute), 29
[suffix](#) ([Projects](#) attribute), 29
[support_year](#) ([Projects](#) attribute), 29

T

[term_counts\(\)](#) (in module [nesta.core.batchables.nlp.vectorizer.run](#)), 26
[terminate_job\(\)](#) ([BatchClient](#) method), 35
[test](#) ([CollectNewTask](#) attribute), 16
[test](#) ([DateTask](#) attribute), 17
[test](#) ([GridTask](#) attribute), 19

[test \(MagSparqlTask attribute\), 18](#)
[test \(QueryMagTask attribute\), 18](#)
[test \(TopicDiscoveryTask attribute\), 21](#)
[timestamp \(Group attribute\), 28](#)
[tokenize_document\(\) \(in module packages.nlp_utils.preprocess\), 8](#)
[topic_perc \(TopicDiscoveryTask attribute\), 21](#)
[TopicDiscoveryTask \(class in nesta.core.routines.meetup.health_tagging.topic_discovery_task\), 21](#)
[topics \(Group attribute\), 27](#)
[total_cost \(Projects attribute\), 29](#)
[total_cost_sub_project \(Projects attribute\), 29](#)
[touch\(\) \(MySQLTarget method\), 37](#)

U

[urlname \(Group attribute\), 27](#)

W

[wait_on_job\(\) \(BatchClient method\), 35](#)