
brownant Documentation

Release 0.1.1

Jiangge Zhang

September 29, 2013

CONTENTS

USER'S GUIDE

1.1 Quick Start

There are some simple crawling applications written with BrownAnt.

1.1.1 The Minimal Demo

This demo could get the download link from the PyPI home page of given project.

```
# example.py
from brownant.app import BrownAnt
from brownant.site import Site
from lxml import html
from requests import Session

site = Site(name="pypi")
http = Session()

@site.route("pypi.python.org", "/pypi/<name>", defaults={"version": None})
@site.route("pypi.python.org", "/pypi/<name>/<version>")
def pypi_info(request, name, version):
    url = request.url.geturl()
    etree = html.fromstring(http.get(url).content)
    download_url = etree.xpath("./div[@id='download-button']/a/@href")[0]

    return {
        "name": name,
        "version": version,
        "download_url": download_url,
    }

app = BrownAnt()
app.mount_site(site)

if __name__ == "__main__":
    from pprint import pprint
    pprint(app.dispatch_url("https://pypi.python.org/pypi/Werkzeug/0.9.4"))
```

And run it, we will get the output:

```
$ python example.py
{'download_url': 'https://.../source/W/Werkzeug/Werkzeug-0.9.4.tar.gz',
 'name': u'Werkzeug',
 'version': u'0.9.4'}
```

1.1.2 The Declarative Demo

With the declarative usage, the workflow will be flexible and readable.

First, we define the “dinergate” in a site supported module:

```
# sites/pypi.py
from brownant.site import Site
from brownant.dinergate import Dinergate
from brownant.pipeline.network import TextResponseProperty
from brownant.pipeline.html import ElementTreeProperty, XPathTextProperty

site = Site(name="pypi")

@site.route("pypi.python.org", "/pypi/<name>/<version>")
class PythonPackageInfo(Dinergate):

    URL_TEMPLATE = "http://pypi.python.org/pypi/{self.name}/{self.version}"

    text_response = TextResponseProperty()
    etree = ElementTreeProperty()
    download_url = XPathTextProperty(
        xpath=".//div[@id='download-button']/a/@href",
        strip_spaces=True, pick_mode="first"
    )

    @property
    def info(self):
        return {"name": self.name, "version": self.version,
                "download_url": self.download_url}
```

And then we define an application instance and mount the site.

```
# app.py
from brownant.app import BrownAnt

app = BrownAnt()
app.mount_site("sites.pypi:site")

if __name__ == "__main__":
    from pprint import pprint
    pkg = app.dispatch_url("https://pypi.python.org/pypi/Werkzeug/0.9.4")
    pprint(pkg.info)
```

And run it, we will get the same output.

API REFERENCE

2.1 Basic API

The basic API included the application framework and routing system (provided by `werkzeug.routing`) of BrownAnt.

2.1.1 brownant.app

`class brownant.app.BrownAnt`

The app which could manage whole crawler system.

`add_url_rule(host, rule_string, endpoint, **options)`

Add a url rule to the app instance.

The url rule is the same with Flask apps and other Werkzeug apps.

Parameters

- **host** – the matched hostname. e.g. “www.python.org”
- **rule_string** – the matched path pattern. e.g. “/news/<int:id>”
- **endpoint** – the endpoint name as a dispatching key such as the qualified name of the object.

`dispatch_url(url_string)`

Dispatch the URL string to the target endpoint function.

Parameters **url_string** – the origin URL string.

Returns the return value of calling dispatched function.

`mount_site(site)`

Mount a supported site to this app instance.

Parameters **site** – the site instance be mounted.

`parse_url(url_string)`

Parse the URL string with the url map of this app instance.

Parameters **url_string** – the origin URL string.

Returns the tuple as (`url, url_adapter, query_args`), the url is parsed by the standard library `urlparse`, the url_adapter is from the werkzeug bound URL map, the query_args is a multidict from the werkzeug.

2.1.2 brownant.request

```
class brownant.request.Request(url, args)
    The crawling request object.
```

Parameters

- **url** (`urllib.parse.ParseResult`) – the raw URL inputted from the dispatching app.
- **args** (`werkzeug.datastructures.MultiDict`) – the query arguments decoded from query string of the URL.

2.1.3 brownant.site

```
class brownant.site.Site(name)
```

The site supported object which could be mounted to app instance.

Parameters **name** – the name of the supported site.

```
play_actions(target)
```

Play record actions on the target object.

Parameters **target** (`BrownAnt`) – the target which receive all record actions, is a brown ant app instance normally.

```
record_action(method_name, *args, **kwargs)
```

Record the method-calling action.

The actions expect to be played on an target object.

Parameters

- **method_name** – the name of called method.
- **args** – the general arguments for calling method.
- **kwargs** – the keyword arguments for calling method.

```
route(host, rule, **options)
```

The decorator to register wrapped function to the brown ant app.

The parameters of this method is compatible with the `add_url_rule()` method.

The registered function or class must be able to be imported by its qualified name.

The right way:

```
@site.route("www.example.com", "/item/<int:item_id>")
def spam(request, item_id):
    pass
```

The wrong way:

```
def egg():
    # the function could not be imported by its qualified name
@site.route("www.example.com", "/item/<int:item_id>")
def spam(request, item_id):
    pass

egg()
```

Parameters

- **host** – the limited host name.
- **rule** – the URL path rule as string.
- **options** – the options to be forwarded to the `werkzeug.routing.Rule` object.

2.1.4 brownant.exceptions

exception `brownant.exceptions.BrownAntException`

The base exception of the brown ant framework.

exception `brownant.exceptions.NotSupported`

Bases: `brownant.exceptions.BrownAntException`

The given URL or other identity is from a platform which not support.

This exception means any url rules of the app which matched the URL could not be found.

2.2 Declarative API

The declarative API is around the “dinergate” and “pipeline property”.

2.2.1 brownant.dinergate

class `brownant.dinergate.Dinergate` (`request`, `http_client=None`, `**kwargs`)

The simple classify crawler.

In order to work with unnamed properties such as the instances of `PipelineProperty`, the meta class `DinergateType` will scan subclasses of this class and name all unnamed members which are instances of `cached_property`.

Parameters

- **request** (`Request`) – the standard parameter passed by app.
- **http_client** (`requests.Session`) – the session instance of python-requests.
- **kwargs** – other arguments from the URL pattern.

URL_TEMPLATE = None

the URL template string for generating crawled target. the `self` could be referenced in the template. (e.g. `“http://www.example.com/items/{self.item_id}?page={self.page}”`)

url

The fetching target URL.

The default behavior of this property is build URL string with the `URL_TEMPLATE`.

The subclasses could override `URL_TEMPLATE` or use a different implementation.

class `brownant.dinergate.DinergateType`

Bases: `type`

The metaclass of `Dinergate` and its subclasses.

This metaclass will give all members are instance of `cached_property` default names. It is because many pipeline properties are subclasses of `cached_property`, but them would not be created by decorating functions. They will has not built-in `__name__`, which may cause them could not cache values as expected.

2.2.2 brownant.pipeline.base

```
class brownant.pipeline.base.PipelineProperty(**kwargs)
    Bases: werkzeug.utils.cached_property
```

The base class of pipeline properties.

There are three kinds of initial parameters.

- The required attribute. If a keyword argument's name was defined in `requiredAttrs`, it will be assigned as an instance attribute.
- The `attr_name`. It is the member of `attr_names`, whose name always end with `_attr`, such as `text_attr`.
- The option. It will be placed at an instance owned `dict` named `options`. The subclasses could set default option value in the `prepare()`.

A workable subclass of `PipelineProperty` should implement the abstract method `provide_value()`, which accept an argument, the instance of `DinerGate`.

Overriding `prepare()` is optional in subclasses.

Parameters `kwargs` – the parameters with the three kinds.

provide_value (`obj`)

The abstract method which should be implemented by subclasses. It provide the value expected by us from the subject instance.

Parameters `obj` (`DinerGate`) – the subject instance.

attr_names = None

the definition of `attr_names`

get_attr (`obj, name`)

Get attribute of the target object with the configured attribute name in the `attr_names` of this instance.

Parameters

- `obj` (`DinerGate`) – the target object.
- `name` – the internal name used in the `attr_names`. (e.g. “`text_attr`”)

options = None

the definition of options

prepare ()

This method will be called after instance initialized. The subclasses could override the implementation.

In general purpose, the implementation of this method should give default value to options and the members of `attr_names`.

Example:

```
def prepare(self):
    self.attr_names.setdefault("text_attr", "text")
    self.options.setdefault("use_proxy", False)
```

requiredAttrs = set([])

the names of required attributes.

2.2.3 brownant.pipeline.network

```
class brownant.pipeline.network.URLQueryProperty(**kwargs)
```

The query argument property. The usage is simple:

```
class MySite(Dinergate):
    item_id = URLQueryProperty(name="item_id", type=int)
```

It equals to this:

```
class MySite(Dinergate):
    @cached_property
    def item_id(self):
        value = self.request.args.get("item_id", type=int)
        if not value:
            raise NotSupported
        return value
```

A failure conversion with given type (`ValueError` be raised) will lead the value fallback to `None`. It is the same with the behavior of the `MultiDict`.

Parameters

- `name` – the query argument name.
- `request_attr` – optional. default: “`request`”.
- `type` – optionl. default: `None`. this value will be passed to `get()`.
- `required` – optionl. default: `True`. while this value be true, the `NotSupported` will be raised for meeting empty value.

```
class brownant.pipeline.network.TextResponseProperty(**kwargs)
```

The text response which returned by fetching network resource.

Getting this property is network I/O operation in the first time. The http request implementations are all provided by `requests`.

The usage example:

```
class MySite(Dinergate):
    foo_http = requests.Session()
    foo_url = "http://example.com"
    foo_text = TextResponseProperty(url_attr="foo_url",
                                    http_client="foo_http",
                                    proxies=PROXIES)
```

Parameters

- `url_attr` – optional. default: “`url`”. it point to the property which could provide the fetched url.
- `http_client_attr` – optional. default: “`http_client`”. it point to an http client property which is instance of `requests.Session`
- `kwargs` – the optional arguments which will be passed to `requests.Session.get()`.

2.2.4 brownant.pipeline.html

```
class brownant.pipeline.html.ElementTreeProperty(**kwargs)
```

The element tree built from a text response property. There is an usage example:

```
class MySite(Dinergate):
    text_response = "<html></html>"
    div_response = "<div></div>"
    etree = ElementTreeProperty()
    div_etree = ElementTreeProperty(text_response_attr="div_response")

    site = MySite(request)
    print(site.etree)  # output: <Element html at 0x1f59350>
    print(site.div_etree)  # output: <Element div at 0x1f594d0>
```

Parameters `text_response_attr` – optional. default: “`text_response`”.

class brownant.pipeline.html.**XPathTextProperty**(**kwargs)
The text extracted from a element tree property by XPath. There is an example for usage:

```
class MySite(Dinergate):
    # omit page_etree
    title = XPathTextProperty(xpath=".//h1[@id='title']/text()", etree_attr="page_etree",
                               strip_spaces=True,
                               pick_mode="first")
    links = XPathTextProperty(xpath=".//*[@id='links']/a/@href", etree_attr="page_etree",
                               strip_spaces=True,
                               pick_mode="join",
                               joiner=" | ")
```

Parameters

- **xpath** – the xpath expression for extracting text.
- **etree_attr** – optional. default: “`etree`”.
- **strip_spaces** – optional. default: *False*. if it be *True*, the spaces in the beginning and the end of texts will be striped.
- **pick_mode** – optional. default: “`join`”, and could be “`join`” or “`first`”. while “`join`” be detected, the texts will be joined to one. otherwise the “`first`” be detected, only the first text would be picked.
- **joiner** – optional. default is a space string. it is no sense in assigning this parameter while the `pick_mode` is not “`join`”. otherwise, the texts will be joined by this string.

INDICES AND TABLES

- *genindex*
- *modindex*
- *search*