Spidermon Documentation

Scrapinghub

Contents

1	Cont	tents	3
	1.1	Installation	
		Getting started	
	1.3	Monitoring your jobs	9
	1.4	Item Validation	14
		Comparing Spider Executions	
	1.6	Settings	19
	1.7	Actions	
	1.8	Release notes	28
Рy	thon]	Module Index	31
Inc	dex		33

Spidermon is a framework to build monitors for Scrapy spiders. It offers the following features:

- It can check the output data produced by Scrapy (or other sources) and verify it against a schema or model that defines the expected structure, data types and value restrictions. It supports data validation based on two external libraries:
 - jsonschema: https://github.com/Julian/jsonschema
 - Schematics: https://github.com/schematics/schematics
- It allows you to define conditions that should trigger an alert based on Scrapy stats.
- It supports notifications via email and Slack.
- It can generate custom reports.

Contents 1

2 Contents

CHAPTER 1

Contents

1.1 Installation

Spidermon's core functionality provides some useful features that allow you to build your monitors on top of it. The library depends on six, jsonschema and python-slugify.

If you want to set up any notifications, additional monitoring dependencies will help with that.

If you want to use schematics validation, you probably want validation.

So the recommended way to install the library is by adding both:

pip install "spidermon[monitoring, validation]"

1.2 Getting started

Spidermon is a monitoring tool for Scrapy spiders that allows you to write monitors that may validate the execution and the results of your spiders.

This tutorial shows how to set up Spidermon to monitor a spider to check if it extracted a minimum number of items and if they follow a defined data model. The results of the spider execution will be sent to a Slack channel.

It is expected that you have a basic knowledge of Scrapy. If that is not the case read the Scrapy Tutorial and come back later. We will also assume that Spidermon is already installed on your system. If that is not the case case, see *Installation*.

1.2.1 Our sample project

You can find the complete code of our tutorial project here.

We are going to scrape quotes.toscrape.com, a website that lists quotes from famous authors. First we need a regular Scrapy project and create a simple spider:

```
$ scrapy startproject tutorial
$ cd tutorial
$ scrapy genspider quotes quotes.toscrape.com
```

And our spider code:

```
# tutorial/spiders/quotes.py
import scrapy
class QuotesSpider(scrapy.Spider):
   name = 'quotes'
    allowed_domains = ['quotes.toscrape.com']
    start_urls = ['http://quotes.toscrape.com/']
    def parse(self, response):
        for quote in response.css('.quote'):
            yield {
                'quote': quote.css('.text::text').get(),
                'author': quote.css('.author::text').get(),
                'author_url': response.urljoin(
                    quote.css('.author a::attr(href)').get()),
                'tags': quote.css('.tag *::text').getall(),
        yield scrapy. Request (
            response.urljoin(
                response.css('.next a::attr(href)').get()
            )
        )
```

1.2.2 Enabling Spidermon

To enable Spidermon in your project, include the following lines in your Scrapy project settings.py file:

```
SPIDERMON_ENABLED = True

EXTENSIONS = {
    'spidermon.contrib.scrapy.extensions.Spidermon': 500,
}
```

1.2.3 Our first monitor

Monitors are similar to test cases with a set of methods that are executed at well defined moments of the spider execution containing your monitoring logic.

Monitors must be grouped into monitor suites which define a list of monitors that need to be executed and the actions to be performed before and after the suite execute all monitors.

Our first monitor will check whether at least 10 items were returned at the end of the spider execution.

Create a new file called *monitors.py* that will contain the definition and configuration of your monitors.

```
# tutorial/monitors.py
from spidermon import Monitor, MonitorSuite, monitors
```

(continued from previous page)

```
@monitors.name('Item count')
class ItemCountMonitor(Monitor):

@monitors.name('Minimum number of items')
def test_minimum_number_of_items(self):
    item_extracted = getattr(
        self.data.stats, 'item_scraped_count', 0)
    minimum_threshold = 10

msg = 'Extracted less than {} items'.format(
        minimum_threshold)
    self.assertTrue(
        item_extracted >= minimum_threshold, msg=msg
    )

class SpiderCloseMonitorSuite(MonitorSuite):
    monitors = [
        ItemCountMonitor,
    ]
```

This suite needs to be executed when the spider closes, so we include it in the SPIDER-MON_SPIDER_CLOSE_MONITORS list in your *settings.py* file:

```
# tutorial/settings.py
SPIDERMON_SPIDER_CLOSE_MONITORS = (
   'tutorial.monitors.SpiderCloseMonitorSuite',
)
```

After executing the spider, you should see the following in your logs:

```
$ scrapy crawl quotes
(...)
INFO: [Spidermon] ----- MONITORS -----
INFO: [Spidermon] Item count/Minimum number of items... OK
INFO: [Spidermon] ------
INFO: [Spidermon] 1 monitor in 0.001s
INFO: [Spidermon] OK
INFO: [Spidermon] ------ FINISHED ACTIONS ------
INFO: [Spidermon] ------
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
INFO: [Spidermon] ------ PASSED ACTIONS ------
INFO: [Spidermon] -----
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
INFO: [Spidermon] ------ FAILED ACTIONS ------
INFO: [Spidermon] -----
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
[scrapy.statscollectors] INFO: Dumping Scrapy stats:
(...)
```

If the condition in your monitor fails, this information will appear in the logs:

```
$ scrapy crawl quotes

(continues on next page)
```

1.2. Getting started 5

(continued from previous page)

```
(...)
INFO: [Spidermon] ----- MONITORS -----
INFO: [Spidermon] Item count/Minimum number of items... FAIL
INFO: [Spidermon] -----
ERROR: [Spidermon]
______
FAIL: Item count/Minimum number of items
Traceback (most recent call last):
 File "/tutorial/monitors.py",
  line 17, in test_minimum_number_of_items
  item_extracted >= minimum_threshold, msg=msg
AssertionError: False is not true : Extracted less than 10 items
INFO: [Spidermon] 1 monitor in 0.001s
INFO: [Spidermon] FAILED (failures=1)
INFO: [Spidermon] ------ FINISHED ACTIONS -----
INFO: [Spidermon] -----
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
INFO: [Spidermon] ------ PASSED ACTIONS ------
INFO: [Spidermon] ------
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
INFO: [Spidermon] ------ FAILED ACTIONS ------
INFO: [Spidermon] ------
INFO: [Spidermon] 0 actions in 0.000s
INFO: [Spidermon] OK
```

1.2.4 Slack notifications

Receiving fail notification in the logs may be helpful during the development but not so useful when you are running a huge number of spiders, so you can define actions to be performed when your spider start or finishes (with or without failures).

Spidermon has some built-in actions but you are free to define your own.

Here we will configure a built-in Spidermon action that sends a pre-defined message to a Slack channel using a bot when a monitor fails.

```
# tutorial/monitors.py
from spidermon.contrib.actions.slack.notifiers import SendSlackMessageSpiderFinished

# (...your monitors code...)

class SpiderCloseMonitorSuite(MonitorSuite):
    monitors = [
        ItemCountMonitor,
    ]

    monitors_failed_actions = [
        SendSlackMessageSpiderFinished,
    ]
```

After enabling the action, you need to provide the Slack credentials in your settings.py file as follows:

```
# tutorial/settings.py
(...)
SPIDERMON_SLACK_SENDER_TOKEN = '<SLACK_SENDER_TOKEN>'
SPIDERMON_SLACK_SENDER_NAME = '<SLACK_SENDER_NAME>'
SPIDERMON_SLACK_RECIPIENTS = ['@yourself', '#yourprojectchannel']
```

If a monitor fails, the recipients provided will receive a message in Slack:

```
16:22 bender APP quotes spider finished with errors! (errors=1)

• Item count/Minimum number of items
```

1.2.5 Item validation

Item validators allows you to match your returned items with predetermined structure ensuring that all fields contains data in the expected format. Spidermon allows you to choose between schematics or JSON Schema to define the structure of your item.

In this tutorial, we will use a schematics model to make sure that all required fields are populated and they are all of the correct format.

First step is to change our actual spider code to use Scrapy items. Create a new file called *items.py*:

```
# tutorial/items.py
import scrapy

class QuoteItem(scrapy.Item):
    quote = scrapy.Field()
    author = scrapy.Field()
    author_url = scrapy.Field()
    tags = scrapy.Field()
```

And then modify the spider code to use the newly defined item:

```
# tutorial/spiders/quotes.py
import scrapy
from tutorial.items import QuoteItem
class QuotesSpider(scrapy.Spider):
   name = 'quotes'
   allowed_domains = ['quotes.toscrape.com']
    start_urls = ['http://quotes.toscrape.com/']
   def parse(self, response):
        for quote in response.css('.quote'):
            item = QuoteItem(
                quote=quote.css('.text::text').get(),
                author=quote.css('.author::text').get(),
                author_url=response.urljoin(
                    quote.css('.author a::attr(href)').get()
                tags=quote.css('.tag *::text').getall()
            yield item
```

8

(continued from previous page)

```
yield scrapy.Request(
    response.urljoin(
        response.css('.next a::attr(href)').get()
    )
)
```

Now we need to create our schematics model in *validators.py* file that will contain all the validation rules:

```
# tutorial/validators.py
from schematics.models import Model
from schematics.types import URLType, StringType, ListType

class QuoteItem(Model):
    quote = StringType(required=True)
    author = StringType(required=True)
    author_url = URLType(required=True)
    tags = ListType(StringType)
```

To allow Spidermon to validate your items, you need to include an item pipeline and inform the name of the model class used for validation:

```
# tutorial/settings.py
ITEM_PIPELINES = {
    'spidermon.contrib.scrapy.pipelines.ItemValidationPipeline': 800,
}
SPIDERMON_VALIDATION_MODELS = (
    'tutorial.validators.QuoteItem',
)
```

After that, every time you run your spider you will have a new set of stats in your spider log providing information about the results of the validations:

```
$ scrapy crawl quotes
(...)
'spidermon/validation/fields': 400,
'spidermon/validation/items': 100,
'spidermon/validation/validators': 1,
'spidermon/validation/validators/item/schematics': True,
[scrapy.core.engine] INFO: Spider closed (finished)
```

You can then create a new monitor that will check these new statistics and raise a failure when we have a item validation error:

```
# monitors.py

from spidermon.contrib.monitors.mixins import StatsMonitorMixin

# (...other monitors...)

@monitors.name('Item validation')
class ItemValidationMonitor(Monitor, StatsMonitorMixin):

@monitors.name('No item validation errors')
def test_no_item_validation_errors(self):
```

(continued from previous page)

You could also set the pipeline to include the validation error as a field in the item (although it may not be necessary, since the validation errors are included in the crawling stats and your monitor can check them once the spiders finishes):

By default, it adds a field called *validation* to the item when the item doesn't match the schema:

```
# tutorial/settings.py
SPIDERMON_VALIDATION_ADD_ERRORS_TO_ITEMS = True
```

The resulted item will look like this:

1.3 Monitoring your jobs

1.3.1 Monitors

Monitors are the main class where you include your monitoring logic. After defining them, you need to include them in a *MonitorSuite*, so they can be executed.

As *spidermon.core.monitors.Monitor* inherits from Python *unittest.TestCase*, you can use all existing assertion methods in your monitors.

In the following example, we define a monitor that will verify whether a minimum number of items were extracted and fails if it is less than the expected threshold.

```
from spidermon import Monitor, monitors
```

(continued from previous page)

```
@monitors.name('Item count')
class ItemCountMonitor(Monitor):

@monitors.name('Minimum items extracted')
def test_minimum_number_of_items_extracted(self):
    minimum_threshold = 100
    item_extracted = getattr(self.data.stats, 'item_scraped_count', 0)
    self.assertFalse(
        item_extracted < minimum_threshold,
        msg='Extracted less than {} items'.format(minimum_threshold)
    )
}</pre>
```

A Monitor instance has the following properties that can be used to help you implement your monitors:

data.stats dict-like object containing the stats of the spider execution data.crawler instance of actual Crawler object

data.spider instance of actual Spider object

1.3.2 Monitor Suites

A *Monitor Suite* groups a set of *Monitor* classes and allows you to specify which actions must be executed at specified moments of the spider execution.

Here is an example of how to configure a new monitor suite in your project:

```
# settings.py
SPIDERMON_SPIDER_OPEN_MONITORS = (
    # list of monitor suites to be executed when the spider starts
)
SPIDERMON_SPIDER_CLOSE_MONITORS = (
    # list of monitor suites to be executed when the spider finishes
)
```

An instance of MonitorSuite defines a set of monitors and actions to be executed after the job finishes its

execution.

name suite name

monitors list of Monitor that will be executed if this suite is enabled.

monitors_finished_actions list of action classes that will be executed when all monitors finished their execution.

monitors_passed_actions list of action classes that will be executed if all monitors passed.

monitors_failed_actions list of action classes that will be executed if at least one of the monitors failed.

order if you have more than one suite enabled in your project, this integer defines the order of execution of the suites

crawler crawler instance

on_monitors_finished(result)

Executed right after the monitors finished their execution and before any other action is executed.

result stats of the spider execution

on_monitors_passed(result)

Executed right after the monitors finished their execution but after the actions defined in *monitors_finished_actions* were executed if all monitors passed.

result stats of the spider execution

on_monitors_failed(result)

Executed right after the monitors finished their execution but after the actions defined in *monitors_finished_actions* were executed if at least one monitor failed.

result stats of the spider execution

1.3.3 The Basic Monitors

Spidermon has some batteries included:)

Check if a job has a expected finish reason.

You can configure the expected reason with the SPIDERMON_EXPECTED_FINISH_REASONS, it should be an iterable of valid finish reasons.

The default value of this settings is: ['finished',].

Check if spider extracted the minimum number of items.

You can configure it using SPIDERMON_MIN_ITEMS setting. There's **NO** default value for this setting, if you try to use this monitor without setting it, it'll raise a NotConfigured exception.

Check for errors in the spider log.

You can configure the expected number of ERROR log messages using SPIDERMON_MAX_ERRORS. The default is 0.

Check for maximum number of unwanted HTTP codes. You can configure it using SPIDERMON_UNWANTED_HTTP_CODES_MAX_COUNT setting or SPIDERMON UNWANTED HTTP CODES setting

This monitor fails if during the spider execution, we receive more than the number of SPIDERMON_UNWANTED_HTTP_CODES_MAX_COUNT setting for at least one of the HTTP Status Codes in the list defined in SPIDERMON_UNWANTED_HTTP_CODES setting.

Default values are:

```
SPIDERMON_UNWANTED_HTTP_CODES_MAX_COUNT = 10
SPIDERMON_UNWANTED_HTTP_CODES = [400, 407, 429, 500, 502, 503, 504, 523, 540, 541]
```

SPIDERMON_UNWANTED_HTTP_CODES can also be a dictionary with the HTTP Status Code as key and the maximum number of accepted responses with that code.

With the following setting, the monitor will fail if more than 100 responses are 404 errors or at least one 500 error:

```
SPIDERMON_UNWANTED_HTTP_CODES = {
   400: 100,
   500: 0,
}
```

1.3.4 Is there a Basic Scrapy Suite ready to use?

Of course, there is! We really want to make it easy for you to monitor your spiders;)

class spidermon.contrib.scrapy.monitors.SpiderCloseMonitorSuite(name=None,

monitors=None, monitors_finished_actions=None, monitors_passed_actions=None, monitors_failed_actions=None, order=None,

crawler=None)

This Monitor Suite implements the following monitors:

- ItemCountMonitor
- ErrorCountMonitor
- FinishReasonMonitor
- UnwantedHTTPCodesMonitor

You can easily enable this monitor after enabling Spidermon:

```
SPIDERMON_SPIDER_CLOSE_MONITORS = (
    'spidermon.contrib.scrapy.monitors.SpiderCloseMonitorSuite',
)
```

1.3.5 Periodic Monitors

Sometimes we don't want to wait until the end of the spider execution to monitor it. For example, we may want to be notified as soon the number of errors reaches a value or close the spider if the time elapsed is greater than expected.

You define your *Monitors* and *Monitor Suites* the same way as before, but you need to provide the time interval (in seconds) between each of the times the *Monitor Suites* is run.

In the following example, we defined a periodic monitor suite that will be executed every minute and will verify if the number of errors found is lesser than a value. If not, the spider will be closed.

First we define a new action that will close the spider when executed:

```
# myproject/actions.py
from spidermon.core.actions import Action

class CloseSpiderAction(Action):

    def run_action(self):
        spider = self.data['spider']
        spider.logger.info("Closing spider")
        spider.crawler.engine.close_spider(spider, 'closed_by_spidermon')
```

Then we create our monitor and monitor suite that verifies the number of errors and then take an action if it fails:

```
# myproject/monitors.py
from myproject.actions import CloseSpiderAction

@monitors.name('Periodic job stats monitor')
class PeriodicJobStatsMonitor(Monitor, StatsMonitorMixin):

@monitors.name('Maximum number of errors reached')
def test_number_of_errors(self):
    accepted_num_errors = 20
    num_errors = self.data.stats.get('log_count/ERROR', 0)

    msg = 'The job has exceeded the maximum number of errors'
    self.assertLessEqual(num_errors, accepted_num_errors, msg=msg)

class PeriodicMonitorSuite(MonitorSuite):
    monitors = [PeriodicJobStatsMonitor]
    monitors_failed_actions = [CloseSpiderAction]
```

The last step is to configure the suite to be executed every 60 seconds:

```
# myproject/settings.py

SPIDERMON_PERIODIC_MONITORS = {
    'myproject.monitors.PeriodicMonitorSuite': 60, # time in seconds
}
```

1.3.6 What to monitor?

These are some of the usual metrics used in the monitors:

- the amount of items extracted by the spider.
- the amount of successful responses received by the spider.

- the amount of failed responses (server-side errors, network errors, proxy errors, etc.).
- the amount of requests that reach the maximum amount of retries and are finally discarded.
- the amount of time it takes to finish the crawl.
- the amount of errors in the log (spider errors, generic errors detected by Scrapy, etc.)
- the amount of bans.
- the job outcome (if it finishes without major issues or if it is closed prematurely because it detects too many bans, for example).
- the amount of items that don't contain a specific field or a set of fields
- the amount of items with validation errors (missing required fields, incorrect format, values that don't match a specific regular expression, strings that are too long/short, numeric values that are too high/low, etc.)

1.4 Item Validation

One useful feature when monitoring a spider is being able to validate your returned items against a defined schema.

Spidermon provides a mechanism that allows you to define an item schema and validation rules that will be executed for each item returned. To enable the item validation feature, the first step is to enable the built-in item pipeline in your project settings:

```
# tutorial/settings.py
ITEM_PIPELINES = {
    'spidermon.contrib.scrapy.pipelines.ItemValidationPipeline': 800,
}
```

After that, you need to choose which validation library will be used. Spidermon accepts schemas defined using schematics or JSON Schema.

1.4.1 With schematics

Schematics is a validation library based on ORM-like models. These models include some common data types and validators, but they can also be extended to define custom validation rules.

Warning: You need to install schematics to use this feature.

```
# Usually placed in validators.py file
from schematics.models import Model
from schematics.types import URLType, StringType, ListType

class QuoteItem(Model):
    quote = StringType(required=True)
    author = StringType(required=True)
    author_url = URLType(required=True)
    tags = ListType(StringType)
```

Check schematics documentation to learn how to define a model and how to extend the built-in data types.

1.4.2 With JSON Schema

JSON Schema is a powerful tool for validating the structure of JSON data. You can define which fields are required, the type assigned to each field, a regular expression to validate the content and much more.

```
Warning: You need to install jsonschema to use this feature.
```

This guide explains the main keywords and how to generate a schema. Here we have an example of a schema for the quotes item from the *tutorial*.

```
"$schema": "http://json-schema.org/draft-07/schema",
"type": "object",
"properties": {
  "quote": {
    "type": "string"
  "author": {
    "type": "string"
  "author_url": {
    "type": "string",
    "pattern": ""
  },
  "tags": {
    "type"
},
"required": [
  "quote",
  "author",
  "author url"
]
```

1.4.3 Settings

These are the settings used for configuring item validation:

SPIDERMON_VALIDATION_ADD_ERRORS_TO_ITEMS

Default: False

When set to True, this adds a field called _validation to the item that contains any validation errors. You can change the name of the field by assigning a name to SPIDERMON_VALIDATION_ERRORS_FIELD:

1.4. Item Validation 15

SPIDERMON_VALIDATION_DROP_ITEMS_WITH_ERRORS

Default: False

Whether to drop items that contain validation errors.

SPIDERMON_VALIDATION_ERRORS_FIELD

Default: _validation

The name of the field added to the item when a validation error happens and SPIDER-MON_VALIDATION_ADD_ERRORS_TO_ITEMS is enabled.

SPIDERMON_VALIDATION_MODELS

Default: None

A list containing the schematics models that contain the definition of the items that need to be validated.

```
# settings.py

SPIDERMON_VALIDATION_MODELS: [
    'myproject.spiders.validators.DummyItemModel'
]
```

If you are working on a spider that produces multiple items types, you can define it as a dict:

```
# settings.py

SPIDERMON_VALIDATION_MODELS: {
    DummyItem: 'myproject.spiders.validators.DummyItemModel',
    OtherItem: 'myproject.spiders.validators.OtherItemModel',
}
```

SPIDERMON_VALIDATION_SCHEMAS

Default: None

A *list* containing the location of the item schema. Could be a local path or a URL.

```
# settings.py

SPIDERMON_VALIDATION_SCHEMAS: [
    '/path/to/schema.json',
    's3://bucket/schema.json',
    'https://example.com/schema.json',
]
```

If you are working on a spider that produces multiple items types, you can define it as a dict:

```
# settings.py

SPIDERMON_VALIDATION_SCHEMAS: {
    DummyItem: '/path/to/dummyitem_schema.json',
    OtherItem: '/path/to/otheritem_schema.json',
}
```

1.4.4 Validation in Monitors

You can build a monitor that checks the validation problems and raises errors if there are too many. You can base it on spidermon.contrib.monitors.mixins.ValidationMonitorMixin which provides methods that can be useful for this. There are 2 groups of methods, for checking all validation errors and specifically for checking missing_required_field errors. All of these methods rely on the job stats, reading spidermon/validation/fields/errors/* entries.

- check_missing_required_fields, check_missing_required_field check that number of missing_required_field errors is less than the specified threshold.
- check_missing_required_fields_percent, check_missing_required_field_percent check that percent of missing_required_field errors is less than the specified threshold.
- check_fields_errors, check_field_errors check that the number of specified (or all) errors is less than the specified threshold.
- check_fields_errors_percent, check_field_errors_percent check that the percent of specified (or all) errors is less than the specified threshold.

All *_field method take a name of one field, while all *_fields method take a list of field names.

Warning: The default behavior for *_fields methods when no field names is passed is to combine error counts for all fields instead of checking each field separately. This is usually not very useful and inconsistent with the behavior when a list of fields is passed, so you should set the correct_field_list_handling monitor attribute to get the correct behavior. This will be the default in some later version.

Some examples:

```
# checks that each of field2 and field3 is missing in no more than 10 items
self.check_missing_required_fields(field_names=['field2', 'field3'], allowed_count=10)
# checks that field2 has errors in no more than 15% of items
self.check_field_errors_percent(field_name='field2', allowed_percent=15)
# checks that no errors is present in any fields
self.check_field_errors_percent()
```

1.5 Comparing Spider Executions

Sometimes it is worthy to compare results from previous executions of your spider. For example, we should be able to track whether the number of items returned is lower than previous executions, which may indicate a problem.

Scrapy allows you to collect stats in the form key/value, but these stats are lost when the spider finishes its execution. Spidermon provides a built-in stats collector that stores the stats of all spider executions in your local file system.

After enabling this stats collector, every spider instance running will have a stats_history attribute containing a list of the stats of the previous spider executions that can be easily accessed in your monitors.

To enable it, include the following code in your project settings:

```
# myproject/settings.py
STATS_CLASS = (
    "spidermon.contrib.stats.statscollectors.LocalStorageStatsHistoryCollector"
)
```

(continued from previous page)

```
# Stores the stats of the last 10 spider execution (default=100)
SPIDERMON_MAX_STORED_STATS = 10
```

The following example shows how we can verify whether the number of items returned in the current spider execution is higher than 90% of the mean of items returned in the previous spider executions.

```
# myproject/monitors.py
from spidermon import Monitor, MonitorSuite, monitors
@monitors.name("History Validation")
class HistoryMonitor(Monitor):
  @monitors.name("Expected number of items extracted")
  def test_expected_number_of_items_extracted(self):
      spider = self.data["spider"]
      total_previous_jobs = len(spider.stats_history)
      if total_previous_jobs == 0:
          return
      previous_item_extracted_mean = (
              previous_job["item_scraped_count"]
              for previous_job in spider.stats_history
          / total_previous_jobs
      items_extracted = self.data.stats["item_scraped_count"]
      # Minimum number of items we expect to be extracted
      minimum_threshold = 0.9 * previous_item_extracted_mean
      self.assertFalse(
          items_extracted <= minimum_threshold,</pre>
          msg="Expected at least {} items extracted.".format(
              minimum_threshold
          ),
class SpiderCloseMonitorSuite (MonitorSuite):
    monitors = [HistoryMonitor]
```

Warning: Running your spider in Scrapy Cloud requires you to manually change some settings in your project:

- 1. Enable DotScrapy Persistence Add-on in your project to keep your .scrapy directory available between job executions.
- 2. STATS_CLASS is overriden by default in Scrapy Cloud. You need to manually include spidermon. contrib.stats.statscollectors.LocalStorageStatsHistoryCollector in your spider settings. The drawback is that you job stats will not be uploaded to Scrapy Cloud interface and will be available only in the job logs.

1.6 Settings

The Spidermon settings allow you to customize the behaviour of your monitors enabling, disabling and configuring features like enabled monitors, monitor actions, item validation and notifications.

1.6.1 Built-in settings reference

Here's a list of all available Spidermons settings, in alphabetical order, along with their default values and the scope where they apply. These settings must be defined in *settings.py* file of your Scrapy project.

SPIDERMON ENABLED

Default: False

Whether to enable Spidermon.

SPIDERMON_EXPRESSIONS_MONITOR_CLASS

SPIDERMON_PERIODIC_MONITORS

```
Default: {}
```

A dict containing the monitor suites that must be executed periodically as key and the time interval (in seconds) between the executions as value.

For example, the following suite will be executed every 30 minutes:

```
SPIDERMON_PERIODIC_MONITORS = {
   'myproject.monitors.PeriodicMonitorSuite': 1800,
}
```

SPIDERMON SPIDER CLOSE MONITORS

List of monitor suites to be executed when the spider closes.

```
SPIDERMON_SPIDER_CLOSE_EXPRESSION_MONITORS
```

SPIDERMON_SPIDER_OPEN_EXPRESSION_MONITORS

```
SPIDERMON SPIDER OPEN MONITORS
```

List of monitor suites to be executed when the spider starts.

SPIDERMON ENGINE STOP MONITORS

List of monitor suites to be executed when the crawler engine is stopped.

SPIDERMON_ENGINE_STOP_EXPRESSION_MONITORS

List of expression monitors to be executed when the crawler engine is stopped.

1.6. Settings

1.7 Actions

By default, when a monitor suite finishes, the pass/fail information is included in the spider logs, which would be enough during development but useless when you are monitoring several spiders.

Spidermon allows you to define actions that are ran after the monitors finish. You can define your own actions or use one of the existing built-in actions.

1.7.1 E-mail action

This action relies on Amazon Simple Email Service to send an e-mail after the monitor suite finishes its execution. In this example, an e-mail will be sent when your monitor suite finishes no matter if it passed or failed:

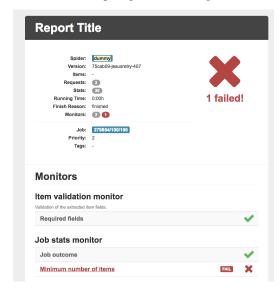
```
from spidermon.contrib.actions.email.ses import SendSESEmail

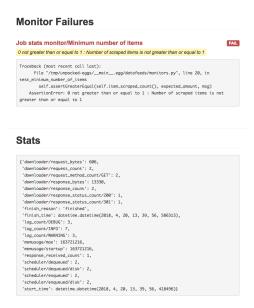
class DummyMonitorSuite(MonitorSuite):
    monitors = [
        DummyMonitor,
    ]

    monitors_finished_actions = [
        SendSESEmail,
    ]
```

By default, Spidermon uses a HTML template that can be altered in *SPIDERMON_BODY_HTML_TEMPLATE* setting. You can use Jinja2 as your template engine.

The result of a report generated using this default template may be seen next:





You can also define actions for when your monitors fails or passes also including actions to the lists monitors_passed_actions and monitors_failed_actions.

The following settings are the minimum needed to make this action works:

SPIDERMON AWS ACCESS KEY

Default: None

AWS Access Key.

SPIDERMON_AWS_SECRET_KEY

AWS Secret Key.

Default: None

SPIDERMON AWS REGION NAME

AWS Region.

Default: us-east-1

SPIDERMON_EMAIL_SENDER

Default: None

Address of the sender of the e-mail notification.

SPIDERMON EMAIL TO

Default: None

List of all recipients of the e-mail notification.

The following settings can be used to improve the action:

1.7. Actions 21

Spidermon Documentation

SPIDERMON BODY HTML

Default: None

SPIDERMON_BODY_HTML_TEMPLATE

String containing the location of the Jinja2 template for the Spidermon email report.

Default reports/email/monitors/result.jinja.

SPIDERMON_BODY_TEXT

SPIDERMON_BODY_TEXT_TEMPLATE

SPIDERMON_EMAIL_BCC

Default: None

SPIDERMON_EMAIL_CONTEXT

Default: None

SPIDERMON_EMAIL_CC

Default: None

SPIDERMON_EMAIL_FAKE

Default: False

If set *True*, the e-mail content will be in the logs but no e-mail will be sent.

SPIDERMON_EMAIL_REPLY_TO

SPIDERMON_EMAIL_SUBJECT

SPIDERMON_EMAIL_SUBJECT_TEMPLATE

1.7.2 Slack action

This action allows you to send custom messages to a Slack channel (or user) using a bot when your monitor suites finish their execution.

To use this action you need to:

1. Install slackclient 1.3 or higher, but lower than 2.0:

```
$ pip install "slackclient>=1.3,<2.0"</pre>
```

Warning: This action **does not** work with slackclient 2.0 or later.

2. Provide the Slack credentials in your settings.py file as follows:

```
# settings.py
SPIDERMON_SLACK_SENDER_TOKEN = '<SLACK_SENDER_TOKEN>'
SPIDERMON_SLACK_SENDER_NAME = '<SLACK_SENDER_NAME>'
SPIDERMON_SLACK_RECIPIENTS = ['@yourself', '#yourprojectchannel']
```

A notification will look like the following one:

```
    bender APP quotes spider finished with errors!
    (errors=1)
    Item count/Minimum number of items
```

The following settings are the minimum needed to make this action works:

```
SPIDERMON SLACK RECIPIENTS
```

List of recipients of the message. It could be a channel or an user.

```
SPIDERMON SLACK SENDER NAME
```

SPIDERMON SLACK SENDER TOKEN

Your Slack token.

Other settings available:

SPIDERMON_SLACK_ATTACHMENTS

SPIDERMON_SLACK_ATTACHMENTS_TEMPLATE

SPIDERMON SLACK FAKE

Default: False

If set True, the Slack message content will be in the logs but nothing will be sent.

SPIDERMON_SLACK_INCLUDE_ATTACHMENTS

SPIDERMON SLACK INCLUDE MESSAGE

SPIDERMON_SLACK_MESSAGE

SPIDERMON SLACK MESSAGE TEMPLATE

SPIDERMON SLACK NOTIFIER INCLUDE ERROR ATTACHMENTS

1.7. Actions 23

SPIDERMON SLACK NOTIFIER INCLUDE OK ATTACHMENTS

SPIDERMON_SLACK_NOTIFIER_INCLUDE_REPORT_LINK

SPIDERMON_SLACK_NOTIFIER_REPORT_INDEX

1.7.3 Job tags action

If you are running your spider using the Scrapy Cloud you are able to include tags in your jobs. Spidermon includes two actions that may be used to add or to remove tags to your jobs depending on the result of the monitoring.

In this example, considering that you defined a *running* tag when you start the job in Scrapy Cloud, if the job passes without errors, it will remove this tag. If the job fails the *failed* tag will be added to the job so you can easily look for failed jobs.

```
# monitors.py
from spidermon.contrib.actions.jobs.tags import AddJobTags, RemoveJobTags

class DummyMonitorSuite(MonitorSuite):
    monitors = [
        DummyMonitor,
    ]

    monitors_passed_actions = [
        RemoveJobTags,
    ]

    monitors_failed_actions = [
        AddJobTags,
    ]
```

```
# settings.py
SPIDERMON_JOB_TAGS_TO_ADD = ['failed', ]
SPIDERMON_JOB_TAGS_TO_REMOVE = ['running', ]
```

By default we have the following settings when using these two actions:

SPIDERMON_JOB_TAGS_TO_ADD

List of tags to be included when AddJobTags is executed.

SPIDERMON_JOB_TAGS_TO_REMOVE

List of tags to be removed when *RemoveJobTags* is executed.

If you want to have different rules adding or removing tags for different results of the monitoring, you need to create a custom action class including the name of the setting that will contain the list of tags that will be included in the job:

```
# monitors.py
from spidermon.contrib.actions.jobs.tags import AddJobTags

class AddJobTagsPassed(AddJobTags):
    tag_settings = 'TAG_TO_ADD_WHEN_PASS'

class AddJobTagsFailed(AddJobTags):
```

(continued from previous page)

```
tag_settings = 'TAG_TO_ADD_WHEN_FAIL'

class DummyMonitorSuite(MonitorSuite):
    monitors = [
         DummyMonitor,
]

monitors_passed_actions = [
         AddJobTagsPassed,
]

monitors_failed_actions = [
         AddJobTagsFailed,
]
```

```
# settings.py
TAG_TO_ADD_WHEN_PASS = ['passed', ]
TAG_TO_ADD_WHEN_FAIL = ['failed', ]
```

1.7.4 File Report action

This action allows to create a file report based on a template. As *E-mail action* you can use Jinja2 as your template engine.

In this example we will create a file called my_report.html when the monitor suite finishes:

```
# settings.py
SPIDERMON_REPORT_TEMPLATE = 'reports/email/monitors/result.jinja'
SPIDERMON_REPORT_CONTEXT = {
    'report_title': 'Spidermon File Report'
}
SPIDERMON_REPORT_FILENAME = 'my_report.html'
```

Settings available:

SPIDERMON_REPORT_CONTEXT

Dictionary containing context variables to be included in your report.

1.7. Actions 25

SPIDERMON REPORT FILENAME

String containing the path of the generated report file.

SPIDERMON REPORT TEMPLATE

String containing the location of the template for the file report.

1.7.5 S3 Report action

This action works exactly like *File Report action* but instead of saving the generated report locally, it uploads it to a S3 Amazon Bucket.

Settings available:

```
SPIDERMON REPORT S3 BUCKET
```

SPIDERMON REPORT S3 CONTENT TYPE

SPIDERMON_REPORT_S3_FILENAME

SPIDERMON_REPORT_S3_MAKE_PUBLIC

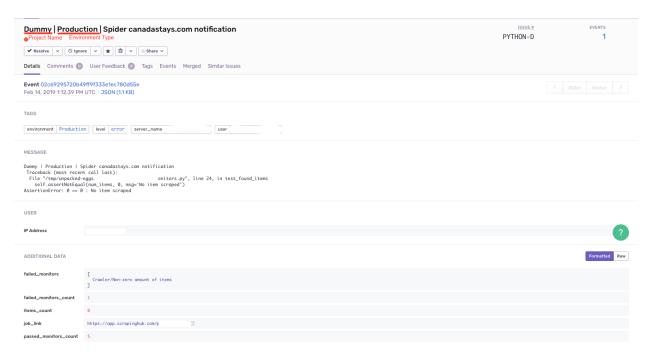
SPIDERMON REPORT S3 REGION ENDPOINT

1.7.6 Sentry action

This action allows you to send custom messages to Sentry when your monitor suites finish their execution. To use this action you need to provide the Sentry DSN in your *settings.py* file as follows:

```
# settings.py
SPIDERMON_SENTRY_DSN = '<SENTRY_DSN_URL>'
SPIDERMON_SENTRY_PROJECT_NAME = '<PROJECT_NAME>'
SPIDERMON_SENTRY_ENVIRONMENT_TYPE = '<ENVIRONMENT_TYPE>'
```

A notification on Sentry will look like the following one:



The following settings are needed to make this action workable:

SPIDERMON_SENTRY_DSN

Data Source Name provided by Sentry, it's a representation of the configuration required by the Sentry SDKs.

SPIDERMON SENTRY PROJECT NAME

Project name to use in notification title.

SPIDERMON SENTRY ENVIRONMENT TYPE

Default: Development

Environment type to use in notification title. It could be set to anything like local, staging, development or production.

SPIDERMON_SENTRY_LOG_LEVEL

Default: error

It could be set to any level provided by Sentry Log Level

SPIDERMON_SENTRY_FAKE

Default: False

If set *True*, the Sentry message will be in the logs but nothing will be sent.

1.7. Actions 27

1.7.7 Custom actions

You can define your own custom actions to be executed by your monitor suites. Just create a class that inherits from *spidermon.core.actions.Action* and implement the *run_action* method.

```
from spidermon.core.actions import Action

class MyCustomAction(Action):
    def run_action(self):
        # Include here the logic of your action
        # (...)
```

1.8 Release notes

1.8.1 1.11.0 (2019-08-02)

- · Allowed per-field checking in ValidationMonitorMixin
- Added option to set AWS Region Name on SES E-Mail action
- Added default Value for 'SPIDERMON_BODY_HTML_TEMPLATE' setting
- Fixed bug in logging of Slack messages when fake setting is enabled
- Enforce lxml 4.3.5 or lower for Python 3.4
- Improved stats history documentation

1.8.2 1.10.2 (2019-07-01)

• Version 1.10.1 with CHANGELOG updated

1.8.3 1.10.1 (2019-07-01)

• Allowed to add absolute location for custom templates

1.8.4 1.10.0 (2019-06-12)

- Added new StatsCollector that access stats data from previous spider executions.
- Added new setting to define the max number of unwanted HTTP status codes allowed in built-in monitor.
- Improved validation error messages with JSON Schema when additional fields are found.
- Made possible to retrieve JSON schema files from external locations.
- Included documentation of periodic monitor suites.
- Fixed bug caused by new slackclient release.
- Other small documentation improvements.

1.8.5 1.9.0 (2019-03-11)

- Add set of built-in basic monitors with the most common test methods to allow start monitoring spiders more straightforward.
- Add SendSentryMessage action to send notifications to Sentry containing the results of Spidermon execution.
- Add SPIDERMON_ENGINE_STOP_MONITORS setting to list monitors to be executed when the Scrapy engine is stopped.
- Fix bug that prevented the use of custom model-level validators in schematics models.
- Refactor JSONSchemaValidator to allow select different versions of JSON Schema.
- · Refactor requirements in setup.py to include missing required dependencies.
- Fix bug caused by backward incompatible change in jsonschema 3.0.0.
- Fix example code of tutorial.
- Install documentation improvements.

1.8.6 1.8.0 (2019-01-08)

- Remove CreateJobReport action.
- Include new documentation and tutorial code.
- Rename internal method in MonitorRunner to fix typo.

1.8.7 1.7.0 (2018-12-04)

- Support universal wheels.
- Skip authentication and recipient settings when running in fake mode.

1.8.8 1.6.0 (2018-11-09)

- Add SPIDERMON_EMAIL_CONTEXT setting to pass custom contexts to email actions.
- Add support for Schematics 2.1.0.

1.8.9 1.5.0 (2018-09-19)

• Convert the job ID tag into a clickable button.

1.8.10 1.4.0 (2018-08-17)

- Avoid requests to get the amount of lines in the log by default, because they consume too much memory and
 they are very slow. You can still use the old behavior adding show_log_count to the context before creating
 the email message.
- Refactor the requirements in setup.py.
- Update the Sphinx configuration.

1.8. Release notes 29

1.8.11 1.3.0 (2018-08-02)

• Add support for periodic monitors in the Scrapy extension.

1.8.12 1.2.0 (2018-04-04)

- Modify ItemValidationPipeline in order to support dict objects in addition to Scrapy.Item objects.
- Refactor ItemValidationPipeline to make it easier to extend this class.

1.8.13 1.1.0 (2018-03-23)

- Add Schematics 2.* support. Note that Schematics 2.0.0 introduced many changes to its API and even some validation rules have a slightly different behaviour in some cases.
- ItemValidationPipeline optimisations for cases where no validators can be applied.

1.8.14 1.0.0 (2018-03-08)

- Add Python 3 support.
- Run tests on Python 2 and Python 3.
- Add dependencies for optional validation features to setup.py.
- Import HubstorageClient from the scrapinghub library if available.
- Replace dash.scrapinghub.com with app.scrapinghub.com.

Backwards Incompatible Changes

- Rename attachements attribute in the SendSlackMessage class to attachments.
- Add the SPIDERMON_ENABLED setting to control if the Scrapy extension should run (note that it is disabled by default).

Python Module Index

S

spidermon.contrib.scrapy.monitors,12

Spidermon I	Documer	ntation
-------------	---------	---------

32 Python Module Index

Index

```
Ε
ErrorCountMonitor
                          (class
                                    in
                                          spider-
        mon.contrib.scrapy.monitors), 11
F
FinishReasonMonitor
                            (class
                                     in
                                          spider-
        mon.contrib.scrapy.monitors), 11
ItemCountMonitor
                         (class
                                   in
                                          spider-
        mon.contrib.scrapy.monitors), 11
M
MonitorSuite (built-in class), 10
on_monitors_failed() (MonitorSuite method), 11
\verb"on_monitors_finished()" (\textit{MonitorSuite method}),
on_monitors_passed() (MonitorSuite method), 11
S
SpiderCloseMonitorSuite (class in spider-
        mon.contrib.scrapy.monitors), 12
spidermon.contrib.scrapy.monitors (mod-
        ule), 11, 12
U
UnwantedHTTPCodesMonitor (class in spider-
        mon.contrib.scrapy.monitors), 11
```