

Scrapy 视图类

在 [Scrapy 源码目录](#) 结构注解中，`website.py` 代表 Web 视图，`webservice.py` 代表 JSON 视图。

视图是使用者与 Scrapy 交互的窗口，本小节我们通过阅读源码来理解 Scrapy 的视图设计。

Web 视图

Scrapy 提供 Web 视图是为了便于使用者监视爬虫运行状态。`website.py` 中有 3 个类，它们分别是：

- 1 1.Root - 根据配置文件初始化项目并设置路由映射关系。
- 2 2.Home - Scrapy 的 Web 首页功能实现及页面渲染。
- 3 3.Jobs - Scrapy 的 jobs 页面功能实现及页面渲染。
- 4

Root 类

Root 类的作用大致为根据配置文件对项目进行初始化并设置路由映射关系，其代码如下：

```
1 class Root(resource.Resource):
2
3     def __init__(self, config, app):
4         resource.Resource.__init__(self)
5         self.debug = config.getboolean('debug', False)
6         self.runner = config.get('runner')
7         logsdir = config.get('logs_dir')
8         itemsdir = config.get('items_dir')
9         local_items = itemsdir and (urlparse(itemsdir).scheme.lower() in ['', 'file'])
10        self.app = app
11        self.nodename = config.get('node_name', socket.gethostname())
12        self.putChild(b'', Home(self, local_items))
13        if logsdir:
14            self.putChild(b'logs', static.File(logsdir.encode('ascii', 'ignore'),
15            'text/plain'))
16        if local_items:
17            self.putChild(b'items', static.File(itemsdir, 'text/plain'))
18            self.putChild(b'jobs', Jobs(self, local_items))
19            services = config.items('services', ())
20            for servName, servClsName in services:
21                servCls = load_object(servClsName)
22                self.putChild(servName.encode('utf-8'), servCls(self))
23            self.update_projects()
24
25        def update_projects(self):
26            self.poller.update_projects()
27            self.scheduler.update_projects()
28
29        @property
```

```

29     def launcher(self):
30         app = IServiceCollection(self.app, self.app)
31         return app.getServiceNamed('launcher')
32
33     @property
34     def scheduler(self):
35         return self.app.getComponent(ISpiderScheduler)
36
37     @property
38     def eggstorage(self):
39         return self.app.getComponent(IEggStorage)
40
41     @property
42     def poller(self):
43         return self.app.getComponent(IPoller)
44

```

代码释义：Root 类继承自 `resource.Resource`，在 `__init__` 方法中通过 `config.get` 读取配置文件中相关配置并赋值给变量，通过 `putChild` 语法将资源与 URL 进行映射（路由），在 `update_projects` 方法中更新项目状态。

Home 类

Home 类的作用大致为 `Scrapyd` 的 `Web` 首页功能实现及页面渲染，其代码大体如下：

```

1  class Home(resource.Resource):
2
3      def __init__(self, root, local_items):
4          resource.Resource.__init__(self)
5          self.root = root
6          self.local_items = local_items
7
8      def render_GET(self, txrequest):
9          vars = {
10              'projects': ', '.join(self.root.scheduler.list_projects())
11          }
12          s = ""
13          <html>
14          <head><title>Scrapyd</title></head>
15          <body>
16          <h1>Scrapyd</h1>
17          <p>Available projects: <b>%(projects)s</b></p>
18          <ul>
19          <li><a href="/jobs">Jobs</a></li>
20          """ % vars
21              if self.local_items:
22                  s += '<li><a href="/items/">Items</a></li>'
23              s += ""
24          ...html ...
25          """ % vars
26          return s.encode('utf-8')
27

```

代码释义: Home 类继承自 `resource.Resource`。在 `render_GET` 方法中, 获取项目列表

`self.root.scheduler.list_projects()`, 然后构造 Web 界面的 HTML 代码。最终将 utf-8 编码的 HTML 代码进行渲染。

Jobs 类

Jobs 类的作用大致为 `Scrapyd` 的 `/jobs` 页面的功能实现及页面渲染, 其代码大体如下:

```
1 class Jobs(resource.Resource):
2
3     def __init__(self, root, local_items):
4         resource.Resource.__init__(self)
5         self.root = root
6         self.local_items = local_items
7
8     cancel_button = ""
9     form...
10    """.format
11
12    header_cols = [
13        'Project', 'Spider',
14        'Job', 'PID',
15        'Start', 'Runtime', 'Finish',
16        'Log', 'Items',
17        'Cancel',
18    ]
19
20    def gen_css(self):
21        css...
22        return '\n'.join(css)
23
24    def prep_row(self, cells):
25        if not isinstance(cells, dict):
26            assert len(cells) == len(self.header_cols)
27        else:
28            cells = [cells.get(k) for k in self.header_cols]
29        cells = ['<td>%s</td>' % ('' if c is None else c) for c in cells]
30        return '<tr>%s</tr>' % ''.join(cells)
31
32    def prep_doc(self):
33        return (
34            html...
35        )
36
37    def prep_table(self):
38        return (
39            html-table...
40        )
41
42    def prep_tab_pending(self):
43        return '\n'.join(
```

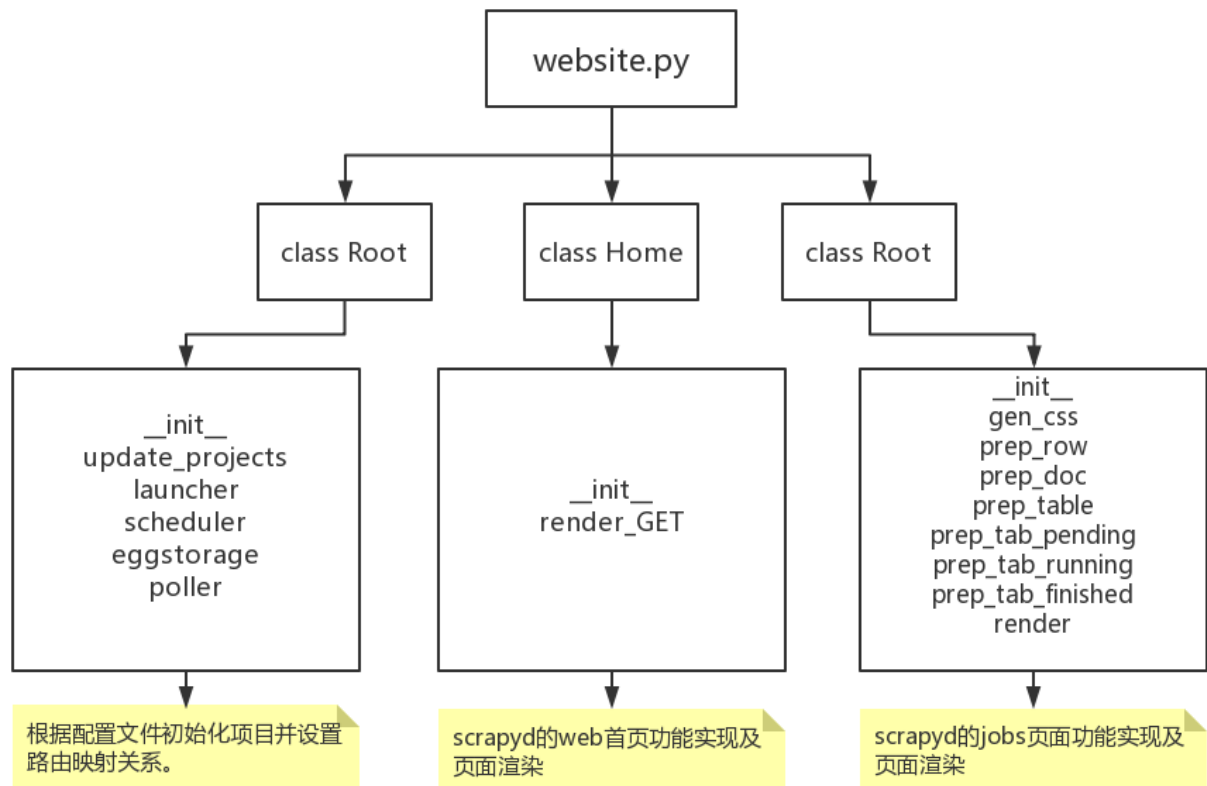
```

44         self.prep_row(dict(
45             Project=project, Spider=m['name'], Job=m['_job'],
46             Cancel=self.cancel_button(project=project, jobid=m['_job']))
47         ))
48         for project, queue in self.root.poller.queues.items()
49         for m in queue.list()
50     )
51
52     def prep_tab_running(self):
53         return '\n'.join(
54             self.prep_row(dict(
55                 ...
56             ))
57             for p in self.root.launcher.processes.values()
58         )
59
60     def prep_tab_finished(self):
61         return '\n'.join(
62             self.prep_row(dict(
63                 ...
64             ))
65             for p in self.root.launcher.finished
66         )
67
68     def render(self, txrequest):
69         doc = self.prep_doc()
70         txrequest.setHeader('Content-Type', 'text/html; charset=utf-8')
71         txrequest.setHeader('Content-Length', len(doc))
72         return doc.encode('utf-8')
73

```

可以看到，Jobs 继承的也是 `resource.Resource`。

Jobs 中的方法是层级调用的关系，从 render 方法开始，调用 `self.prep_doc()` 以获取 HTML 文本，而 `prep_doc` 方法则通过调用 `gen_css` 和 `prep_table` 以获取 CSS 样式及不同状态的爬虫数与爬虫运行记录。所以我们可以绘制出 `website.py` 文件中各个类与方法的组织结构图：



resource.Resource

`resource.Resource` 是 Home 和 Jobs 的父类，它的代码是怎样的？

在 PyCharm 中，通过 `Ctrl/Command+鼠标左键` 可以跟进代码。我们跟进 Resource，代码如下：

```

1  @implementer(IResource)
2  class Resource:
3      entityType = IResource
4
5      server = None
6
7      def __init__(self):
8          """
9          Initialize.
10         """
11         self.children = {}
12
13         isLeaf = 0
14
15         ### Abstract Collection Interface
16
17         def listStaticNames(self):
18             return list(self.children.keys())
19
20         def listStaticEntities(self):
21             return list(self.children.items())
  
```

```
22
23     def listNames(self):
24         return list(self.listStaticNames()) + self.listDynamicNames()
25
26     def listEntities(self):
27         return list(self.listStaticEntities()) + self.listDynamicEntities()
28
29     def listDynamicNames(self):
30         return []
31
32     def listDynamicEntities(self, request=None):
33         return []
34
35     def getStaticEntity(self, name):
36         return self.children.get(name)
37
38     def getDynamicEntity(self, name, request):
39         if name not in self.children:
40             return self.getChild(name, request)
41         else:
42             return None
43
44     def delEntity(self, name):
45         del self.children[name]
46
47     def reallyPutEntity(self, name, entity):
48         self.children[name] = entity
49
50     # Concrete HTTP interface
51
52     def getChild(self, path, request):
53
54         return NoResource("No such child resource.")
55
56
57     def getChildWithDefault(self, path, request):
58
59         if path in self.children:
60             return self.children[path]
61         return self.getChild(path, request)
62
63
64     def getChildForRequest(self, request):
65         warnings.warn("Please use module level getChildForRequest.", DeprecationWarning, 2)
66         return getChildForRequest(self, request)
67
68
69     def putChild(self, path, child):
70
71         self.children[path] = child
72         child.server = self.server
73
74
```

```

75     def render(self, request):
76
77         m = getattr(self, 'render_' + nativeString(request.method), None)
78         if not m:
79             try:
80                 allowedMethods = self.allowedMethods
81             except AttributeError:
82                 allowedMethods = _computeAllowedMethods(self)
83             raise UnsupportedMethod(allowedMethods)
84         return m(request)
85
86
87     def render_HEAD(self, request):
88
89         return self.render_GET(request)
90
91

```

其中我们需要关注的方法有以下几个：

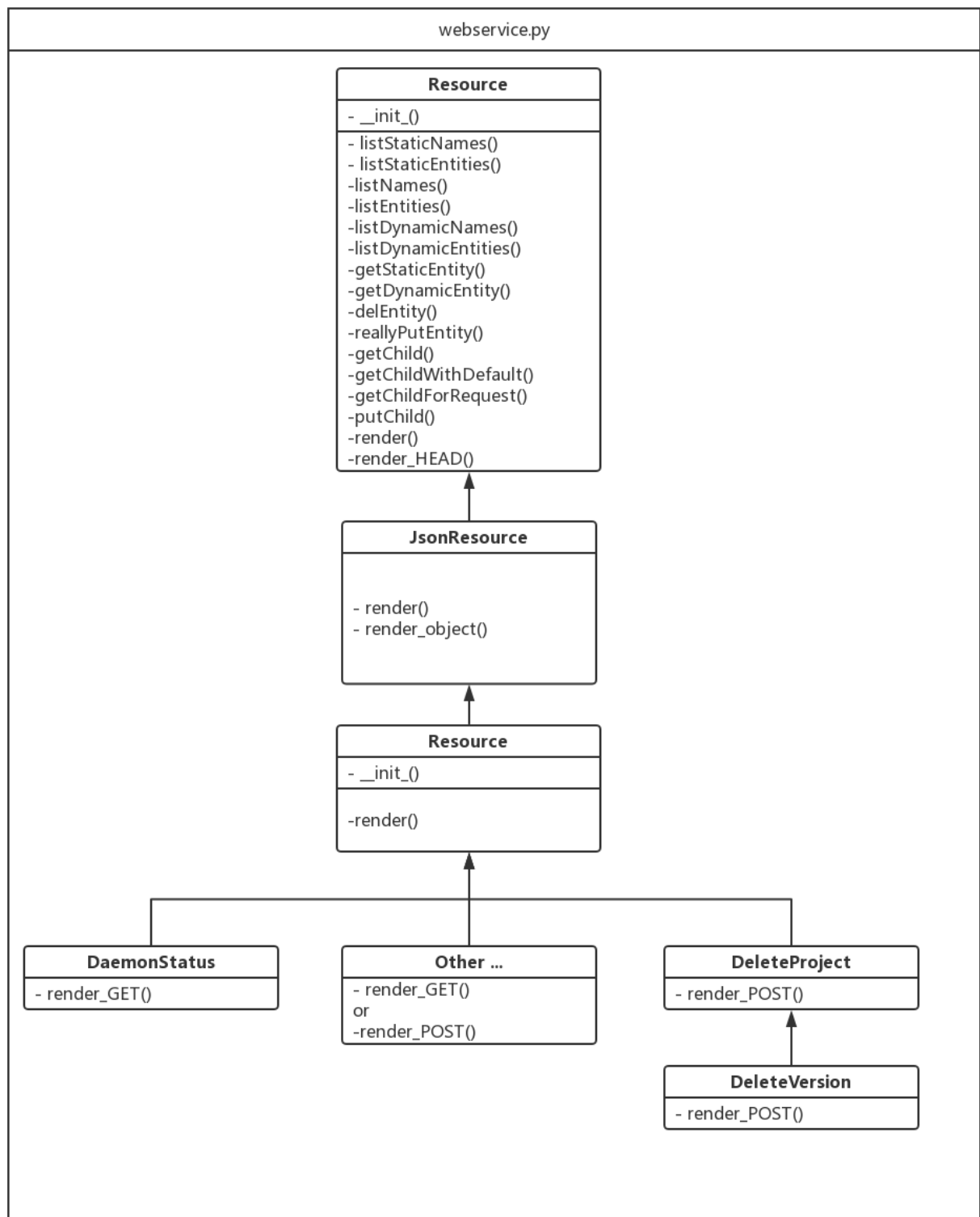
```

1  * putchild - twisted 语法，通过此方法设置路由。
2  * render - 负责页面渲染。
3  * render_HEAD - 将请求转发给 `render_GET`
4

```

JSON 视图

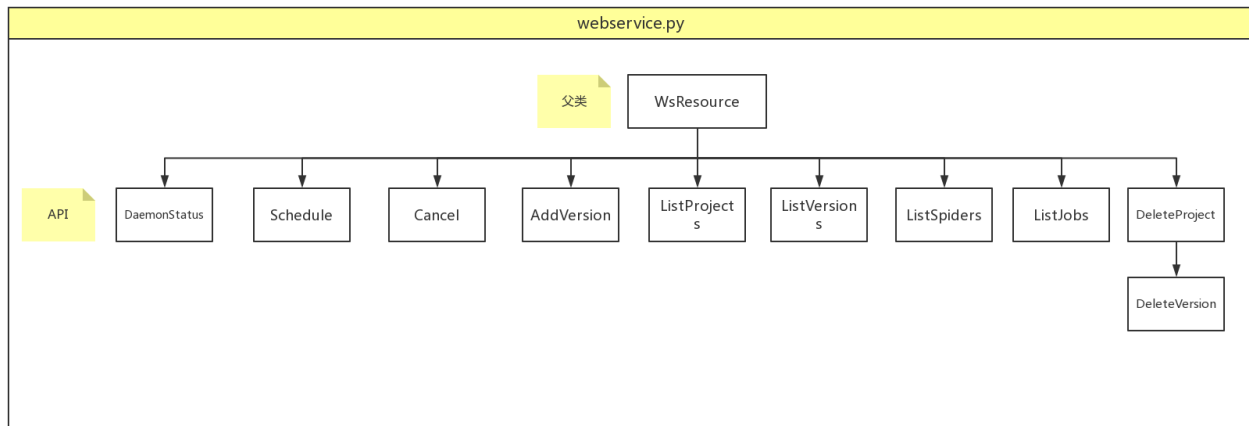
Scrapyd 提供了很多 API，通过 API 的返回信息能够知晓爬虫的运行状态或结果等信息，`webservice.py` 中视图类与父类的继承关系 UML 类图如下：



以下是各个类的大致功能释义：


```
1  * WsResource - 继承自 JsonResponse, 是大多数 API 的父类, 它默认返回 json 数据。
2  * DaemonStatus - 对应 API 中的 daemonstatus.json, 用以检查 Scrapyd 服务的爬虫状态及对应数量。
3  * Schedule - 对应 API 中的 schedule.json, 用以启动指定的爬虫。
4  * Cancel - 对应 API 中的 cancel.json, 用以取消指定的爬虫。
5  * AddVersion - 对应 API 中的 addversion.json, 用以为项目添加版本, 如果项目不存在则创建项目。
6  * ListProjects - 对应 API 中的 listprojects.json, 用以查看当前已部署的项目名称。
7  * ListVersions - 对应 API 中的 listversions.json, 用以查看指定项目的版本号。
8  * ListSpiders - 对应 API 中的 listspiders.json, 用以查看指定项目的爬虫名称。
9  * ListJobs - 对应 API 中的 listjobs.json, 用以查看指定项目下爬虫的运行状态。
10 * DeleteProject - 对应 API 中的 delversion.json, 用以删除指定项目的指定版本, 版本不存在则删除项目。
11 * DeleteVersion - 对应 API 中的 delproject.json, 用以删除指定项目及其已上传的所有版本。
12
```

以及对应的类组织结构图:



可以看到, 大多数 API 类继承的是 `WsResource`, 并且通过重写 `render` 方法将数据渲染。

在 PyCharm 中, 通过 `Ctrl/Command+鼠标左键` 可以跟进代码, 我们跟进 `WsResource`:

```

22     except Exception as e:
23         if self.root.debug:
24             return traceback.format_exc().encode('utf-8')
25         log.err()
26         r = {"node_name": self.root.nodename, "status": "error", "message": str(e)}
27         return self.render_object(r, txrequest).encode('utf-8')
28
29
30 class DaemonStatus(WsResource):
31
32     def render_GET(self, txrequest):
33         pending = sum(q.count() for q in self.root.poller.queues.values())
34         running = len(self.root.launcher.processes)
35         finished = len(self.root.launcher.finished)
36
37         return {"node_name": self.root.nodename, "status": "ok", "pending": pending, "running":
38
39
40 class Schedule(WsResource):
41
42     def render_POST(self, txrequest):
43         args = native_stringify_dict(copy(txrequest.args), keys_only=False)
44         settings = args.pop('setting', [])
45         settings = dict(x.split('=', 1) for x in settings)

```

代码如下：

```

1 class WsResource(JsonResource):
2
3     def __init__(self, root):
4         JsonResource.__init__(self)
5         self.root = root
6
7     def render(self, txrequest):
8         try:
9             return JsonResource.render(self, txrequest).encode('utf-8')
10        except Exception as e:
11            if self.root.debug:
12                return traceback.format_exc().encode('utf-8')
13            log.err()
14            r = {"node_name": self.root.nodename, "status": "error", "message": str(e)}
15            return self.render_object(r, txrequest).encode('utf-8')
16

```

它继承自 `JsonResource` 并且通过 `JsonResource` 的 `render` 返回 `json` 格式数据，如果报错则返回报错信息，而跟进 `JsonResource` 发现它其实也是继承自 `resource.Resource` 类。所以这里可以得出几个结论：

- `Resource` 是 Scrapy 中视图部分最底层的类，且其 `render` 默认返回 `bytes` 类型数据。
- `JsonResource` 以及 `WsResource` 都是通过重写 `render` 以实现返回 `JSON` 类型数据。
- 或许可以通过自定义视图类，继承自 `Resource`，也通过重写 `render` 方法实现兼容 `bytes` 与 `JSON` 的视图类。