

# 官方源码剖析与 API 详解

本小节我们将对 Scrapyd 原有的 API 做个详细的剖析，只有在了解它原有设计思想与代码风格后，我们才能够照猫画虎，设计出风格相似的代码模块，为之后我们自定义 API 和编写权限验证功能打下基础。

API 相关代码在 `webservice.py` 文件中。

根据前面的小节，我们知道 Scrapyd 的视图分为 HTML 和 JSON 两种。我们所看到的数据呈现都是由视图类处理的，比如：

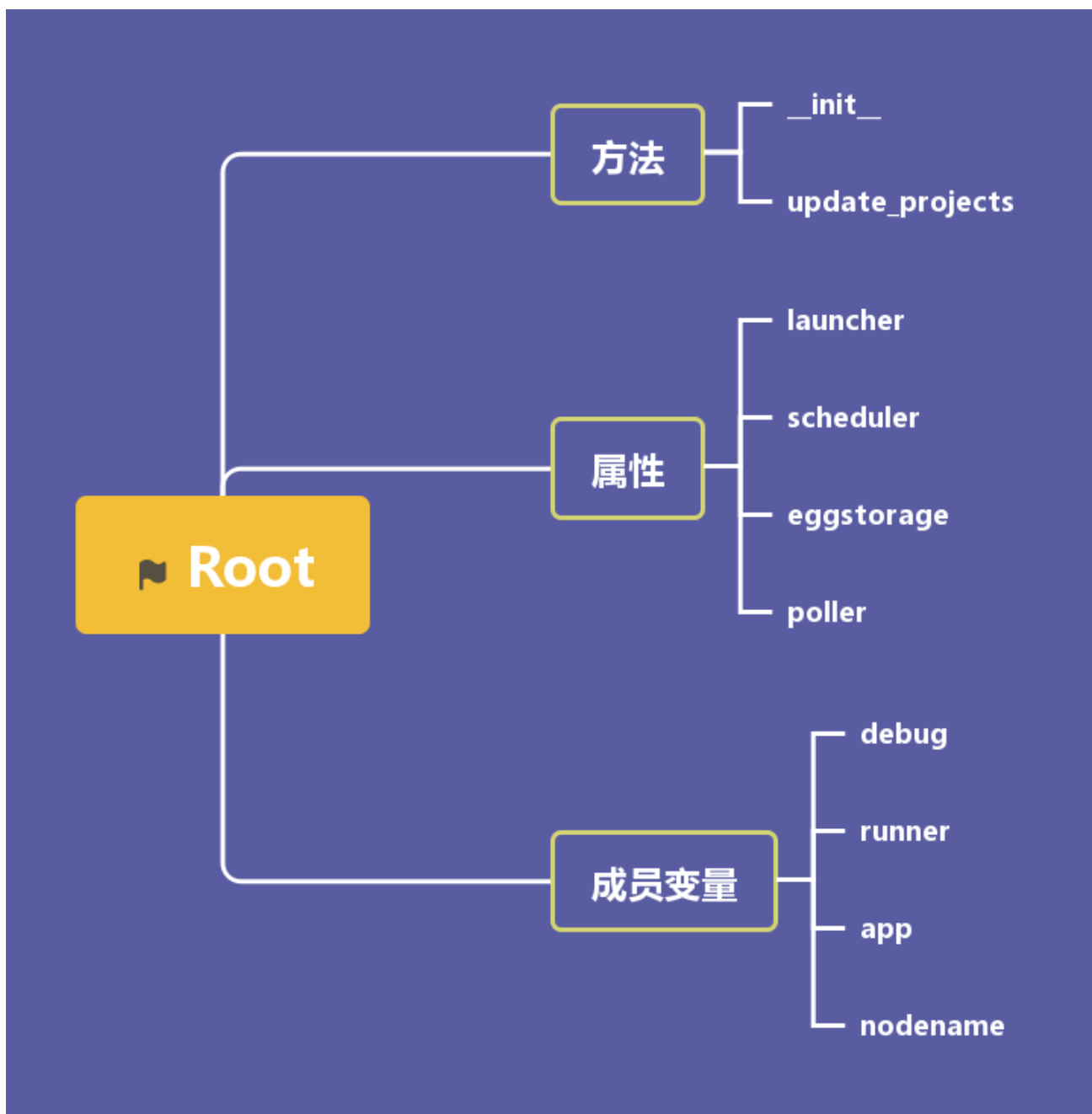
```
1 * 当我们访问根目录的时候，对应的的是 Home 这个类；
2 * 而访问 Jobs 的时候，对应的的是 Jobs 这个类；
3 * 而爬虫的 API 则是，如启动爬虫的 Schedule 类和查看爬虫列表的 ListSpiders 类。
4
```

## HTML 视图

HTML 视图类都写在 Scrapyd 目录下的 `website.py` 文件中，里面有三个类：`Root`、`Home`、`Jobs`。

### Root 类

Root 完成了 `web` 路由设置，Scrapyd 一些基础配置的读取设置，日志目录和 `Item` 目录与路由配置、项目信息更新等任务，是 Scrapyd 的重要组成部分，其代码结构如下图所示：



## Web 路由设置

在 Root 类的 `__init__` 方法中 Home 类以及 Jobs 类的路由是通过 Twisted 的 `putChild` 进行配置的：

```
1 self.putChild(b'jobs', Jobs(self, local_items))
2 self.putChild(b'', Home(self, local_items))
3
```

## 日志与 Item 目录配置

同样在 `__init__` 方法中，先读取日志与 Item 的目录：

```
1 logsdir = config.get('logs_dir')
2 itemsdir = config.get('items_dir')
3
```

然后为它们设置路由：

```
1 self.putChild(b'items', static.File(itemsdir, 'text/plain'))
2 self.putChild(b'logs', static.File(logsdir.encode('ascii', 'ignore'), 'text/plain'))
3
```

## 项目信息更新

当项目变动，比如增删 projects，就会通过 `update_projects()` 方法进行更新：

```
1 def update_projects(self):
2     self.poller.update_projects()
3     self.scheduler.update_projects()
4 @property
5 def poller(self):
6     return self.app.getComponent(IPoller)
7
8 @property
9 def scheduler(self):
10    return self.app.getComponent(ISpiderScheduler)
11
12
```

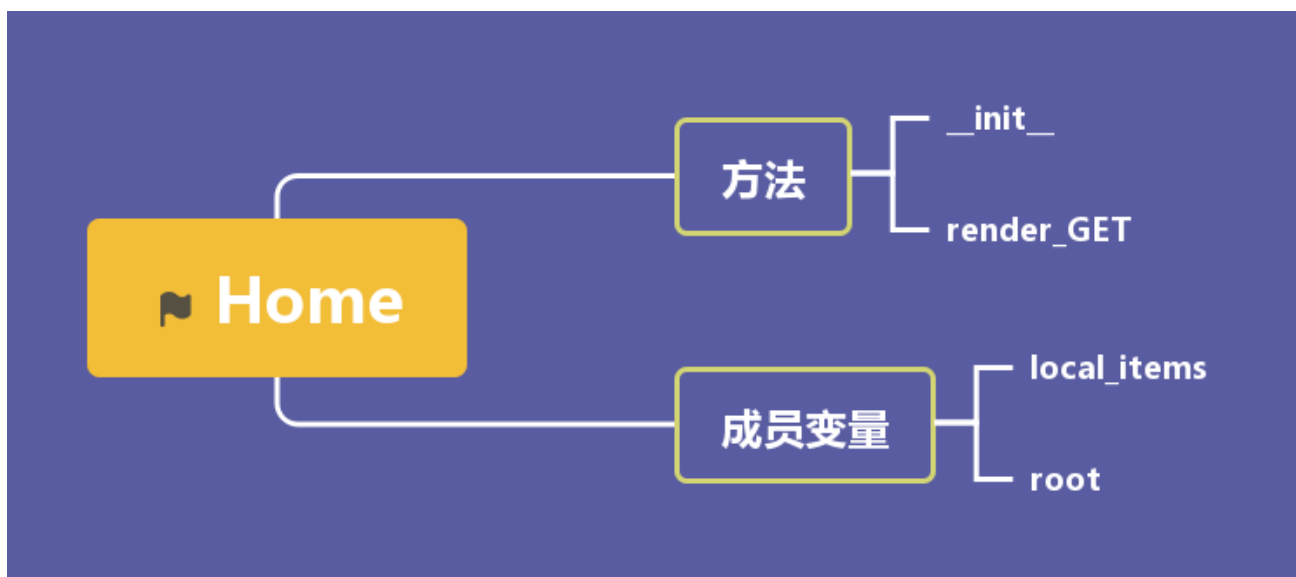
## Scrapyd 的其他一些基础配置等

比如读取配置文件并且将其赋值给变量：

```
1 self.debug = config.getboolean('debug', False)
2 self.runner = config.get('runner')
3 services = config.items('services', ())
4
```

## Home 类

Home 负责 Scrapyd 首页的呈现，当我们访问 `http://localhost:6800` 时看到的界面，就是访问 Home 类，它完成了页面 HTML 布局以及当前已有项目 `projects` 的名称列表展示。



## `__init__` 方法

Home 类继承自 `resource.Resource`，并且重写了 `__init__` 方法，以定义一个 Web 可访问资源：

```
1 def __init__(self, root, local_items):
2     resource.Resource.__init__(self)
3     self.root = root
4     self.local_items = local_items
5
```

## 对于页面呈现的建议

Resource 中 `render` 对于页面呈现的建议：

```
1 """
2 I delegate to methods of self with the form 'render_METHOD'
3 where METHOD is the HTTP that was used to make the
4 request. Examples: render_GET, render_HEAD, render_POST, and
5 so on. Generally you should implement those methods instead of
6 overriding this one.
7 """
8
```

## HTML 布局的实现以及已有项目列表展示

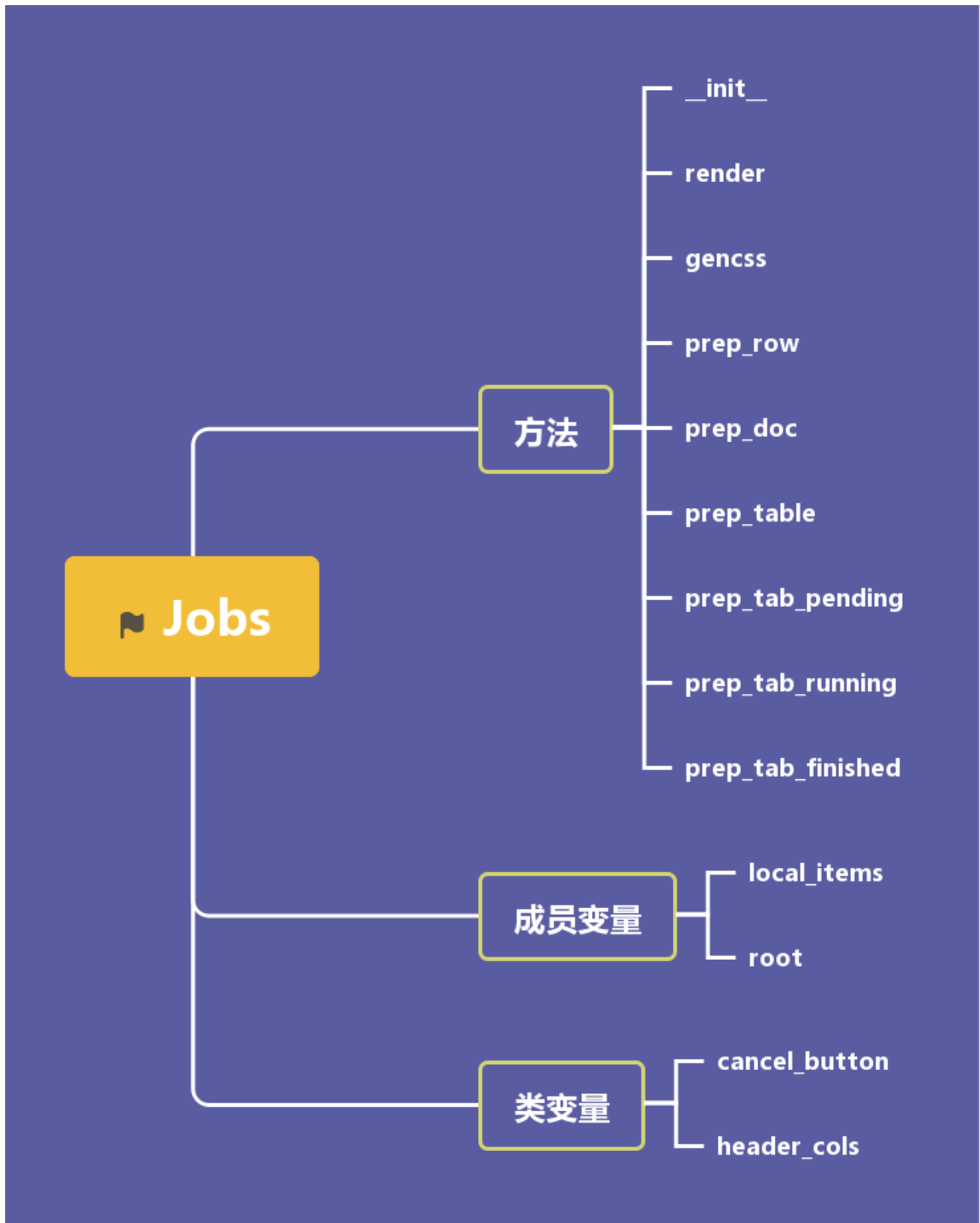
使用 `render` 建议的方式，用 `render_GET` 来完成页面的呈现：

```
1 def render_GET(self, txrequest):
2     vars = {'projects': ', '.join(self.root.scheduler.list_projects())}
3     s = "....."
4     return s.encode('utf-8')
5
```

## Jobs 类

Jobs 负责 Scrapyd 首页爬虫运行状态展示、日志记录以及取消爬虫运行，当我们访问

`http://localhost:6800/jobs/` 时看到的界面，就是 Jobs 类。同样，它也是继承了 `resource.Resource`。



Cancel 与 Header

它设定取消按钮以及爬虫运行信息的列名：

```
1 cancel_button = ""
2 <form method="post" action="/cancel.json">
3     .....
4     .....
5     """.format
6
7 header_cols = [
8     'Project', 'Spider', 'Job', 'PID', 'Start',
9     'Runtime', 'Finish', 'Log', 'Items', 'Cancel',
10 ]
11
```

## 设定 CSS 样式

为爬虫日志表格设定 CSS 样式：

```
1 def gen_css(self):
2     css = [
3         '#jobs>thead td {text-align: center; font-weight: bold}',
4         '#jobs>tbody>tr:first-child {background-color: #eee}'
5     ]
6     return '\n'.join(css)
7
```

## 生成表头

生成爬虫日志表格的表头数据：

```
1 def prep_row(self, cells):
2     if not isinstance(cells, dict):
3         assert len(cells) == len(self.header_cols)
4     else:
5         cells = [cells.get(k) for k in self.header_cols]
6     cells = ['<td>%s</td>' % ('' if c is None else c) for c in cells]
7     return '<tr>%s</tr>' % ''.join(cells)
8
```

## Pending 状态爬虫列表

```
1 def prep_tab_pending(self):
2     return '\n'.join(
3         .....
4     )
5
```

## 正在运行的爬虫列表

```

1 def prep_tab_running(self):
2     return '\n'.join(
3         .....
4         for p in self.root.launcher.processes.values()
5     )
6

```

## 运行完毕的爬虫列表及日志记录

```

1 def prep_tab_finished(self):
2     return '\n'.join(
3         .....
4         for p in self.root.launcher.finished
5     )
6

```

## 生成爬虫运行信息表

调用以上功能，生成爬虫运行信息表格：

```

1 def prep_table(self):
2     return (
3         '<table id="jobs" border="1">'
4         self.prep_row(self.header_cols)
5         self.prep_tab_pending()
6         self.prep_tab_running()
7         self.prep_tab_finished()
8         '..... ..'
9     )
10

```

## 页面渲染

使用 render 方法，将生成的爬虫运行信息呈现到 `localhost:6800/jobs/` 页面：

```

1 def render(self, txrequest):
2     doc = self.prep_doc()
3     txrequest.setHeader('Content-Type', 'text/html; charset=utf-8')
4     txrequest.setHeader('Content-Length', len(doc))
5     return doc.encode('utf-8')
6

```

## JSON 视图

JSON 视图类都写在 Scrapyd 目录下的 `webservice.py` 文件中，里面除了官方文档中提到的 API 外，还有它们父类 `WsResource`。这里我挑选 3 个 API 对应的类进行讲解。

### ListProjects 类

ListProjects 类的作用是返回当前 Scrapy 服务器上的爬虫项目列表。

```
1 class ListProjects(WsResource):
2
3     def render_GET(self, txrequest):
4         projects = list(self.root.scheduler.list_projects())
5         return {"node_name": self.root.nodename, "status": "ok", "projects": projects}
6
7
```

它使用 `render_GET` 方法，所以在浏览器可以输入 `http://localhost:6800/listprojects.json` 来访问并获得 JSON 格式的结果，比如：

```
1 {"status": "ok", "projects": ["AlibabaProject", "JDProject"]}
2
```

## ListSpiders 类

ListSpiders 的作用是返回当前 Scrapy 服务器上指定项目名的爬虫列表。

```
1 class ListSpiders(WsResource):
2
3     def render_GET(self, txrequest):
4         args = native_stringify_dict(copy(txrequest.args), keys_only=False)
5         project = args['project'][0]
6         version = args.get('_version', [''])[0]
7         spiders = get_spider_list(project, runner=self.root.runner, version=version)
8         return {"node_name": self.root.nodename, "status": "ok", "spiders": spiders}
9
```

它也是使用 `render_GET` 方法，但是它需要携带项目名称作为请求参数。如：

`http://localhost:6800/listspiders.json?project=AlibabaProject` 返回的结果同样是 json 格式，通过返回的 json 数据，我们知道 AlibabaProject 项目中当前有哪些爬虫。

```
1 {"status": "ok", "spiders": ["taobao", "1688", "tmall"]}
2
```

## DeleteProject 类

DeleteProject 的作用是用于删除 Scrapy 上已有的爬虫项目。



```

1 class DeleteProject(WsResource):
2
3     def render_POST(self, txrequest):
4         args = native_stringify_dict(copy(txrequest.args), keys_only=False)
5         project = args['project'][0]
6         self._delete_version(project)
7         UtilsCache.invalidate_cache(project)
8         return {"node_name": self.root.nodename, "status": "ok"}
9
10    def _delete_version(self, project, version=None):
11        self.root.eggstorage.delete(project, version)
12        self.root.update_projects()
13

```

它使用的则是 `render_POST` 方法，所以请求的时候我们必须以 post 方式进行，并且它要求携带项目名称作为参数值，如：`http://localhost:6800/delproject.json -d project=AlibabaProject` 在成功删除项目后，更新爬虫项目列表，

```

1 self.root.update_projects()
2

```

如果操作全部成功，我们得到的响应为：

```

1 {"status": "ok"}
2

```

## 视图父类与 Resource

JSON 视图和 HTML 视图的父类是不同的。在 HTML 视图部分，Home 和 Jobs 都继承自 `resource.Resource`，而 JSON 视图部分则继承自 `WsResource`。

### WsResource

```

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
16            return self.render_object(r, txrequest).encode('utf-8')

```

它定义了 render 方法，默认返回的是 json 格式，当返回 json 格式出错时返回报错信息。`JsonResource`中为 json 格式定义了一些头信息，并且将 return 的数据转为 json 格式数据。

## resource.Resource

`Resource` 中有很多方法，它的主要功能是定义一个可访问的 Web 资源，并且提供 HTTP 请求的标准、URL 路由设定标准，如 `putChild` 等，下图为 `Resource` 类的结构图：

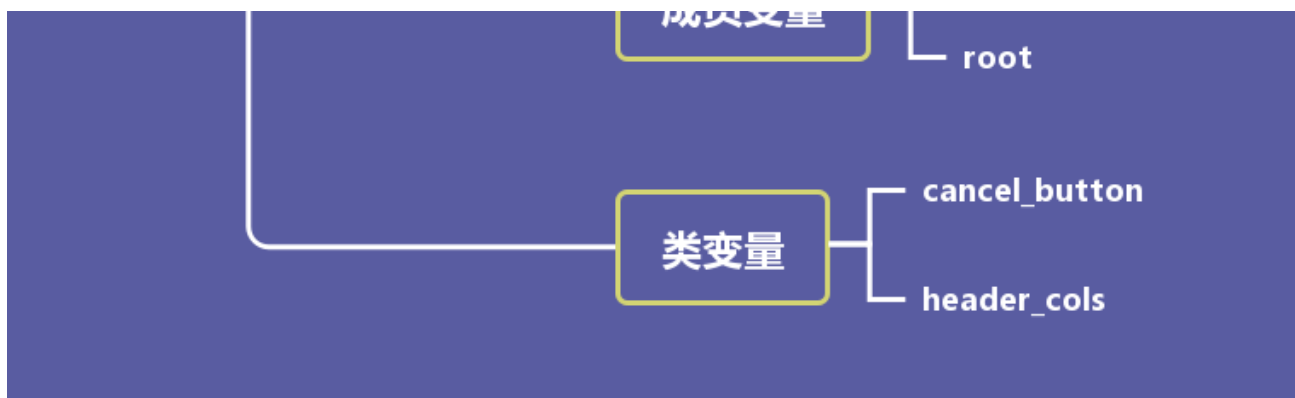
# Jobs

## 方法

- `_init_`
- `render`
- `listStaticNames`
- `listStaticEntities`
- `listNames`
- `listEntities`
- `listDynamicNames`
- `listDynamicEntities`
- `getStaticEntity`
- `getDynamicEntity`
- `delEntity`
- `reallyPutEntity`
- `getChild`
- `getChildWithDefault`
- `getChildForRequest`
- `putChild`

## 成员变量

- `local_items`



## 小结

本小节通过阅读 Scrapyd 中负责 HTML 视图的 `Root`、`Home`、`Jobs` 类和阅读 `API` 中的几个类，知道了 Scrapyd 的 视图及 `API` 构成，并通过阅读其 父类 加深了对 Scrapyd 视图的理解。