

---

# **AyugeSpiderTools**

*Release 3.0.1*

**shengchenyang**

**May 18, 2023**



# 第一步

<b>1</b>	<b>Getting help</b>	<b>3</b>
<b>2</b>	<b>第一步</b>	<b>5</b>
2.1	AyugeSpiderTools 一目了然 . . . . .	5
2.2	安装指南 . . . . .	10
2.3	AyugeSpiderTools 教程 . . . . .	11
2.4	例子 . . . . .	16
<b>3</b>	<b>基本概念</b>	<b>19</b>
3.1	命令行工具 . . . . .	19
3.2	Item . . . . .	23
3.3	Item Loaders . . . . .	28
<b>4</b>	<b>内置服务</b>	<b>31</b>
4.1	Logging . . . . .	31
<b>5</b>	<b>扩展 scrapy</b>	<b>35</b>
5.1	downloader-middleware . . . . .	35
5.2	pipelines . . . . .	41
<b>6</b>	<b>构建你的专属库</b>	<b>45</b>
6.1	How-To-Build-Your-Own-Library . . . . .	45
<b>7</b>	<b>补充说明</b>	<b>47</b>
7.1	Release notes . . . . .	47



在此之前，我们需要了解 Scrapy 是一个快速的高级 [网络爬虫](#) 和 [网页抓取](#) 的框架，用于抓取网站并从其网页中提取结构化数据。它可以用于广泛的目的，从数据挖掘到监控和自动化测试。

AyugeSpiderTools 是充分发挥 Scrapy 的模板功能的一个工具库，可以很方便地管理 Scrapy 项目，比如可以使得我们方便地生成 Scrapy 项目结构，当使用本库内置工具时可以不用每次手动创建 items, middlewares, pipelines, settings 等，内置了比较通用和常见的 middlewares 和 pipelines。但如果你常用的功能不在此库中，你可以自行添加修改后 build 成为你专属的工具库。



## GETTING HELP

遇到麻烦？请优先尝试使用以下方式提问！

- 请在本库 [ayugespidertools github](#) 上提 `issues`。
- 除非一些功能性 `bug`，其它的功能依赖于 `scrapy`，你或许能在 [scrapy issues](#) 或社区中找到答案。
- 若有其它问题也可尝试 [邮箱](#) 联系。





## 2.1 AyugeSpiderTools 一目了然

AyugeSpiderTools 是 Scrapy 的功能扩展模块，对其 spider, item, middleware, pipeline 等模块中的常用功能进行模板化生成和配置。比如生成常见的 spider，运行 sh 和 settings 配置等脚本和固定项目文件结构；也对其不同模块进行功能扩展，比如给 spider 挂上 Mysql engine 的单例句柄可用于 yield 入库前的去重方式之一，给 pipeline 添加自动生成 Mysql 存储场景下所依赖的数据库、数据表、数据字段及注释，也可以解决常见的（字段编码，Data too long，存储字段不存在等等）错误场景。还有很多功能，请在其 Github 上查看。

AyugeSpiderTools 相关信息：

1. 具体请查看对应链接：[\[AyugeSpiderTools\] \(https://github.com/shengchenyang/AyugeSpiderTools\)](https://github.com/shengchenyang/AyugeSpiderTools)

### 2.1.1 注意：

如果你觉得某些功能实现未达到你的期望，比如某些中间件或管道等的实现方法你有更好的方式，你完全可以自行修改和 build，让其成为你个人或小组中的专属库。你可以修改任何你觉得有必要的部分，包括库名在内，希望本库能给你在爬虫开发或 scrapy 扩展开发方面有所指引。

当然，你也可以选择给此项目做出贡献，比如增加或优化某些功能等，但在此之前请提相关的 ISSUES 经确认后开发并提交 PULL REQUESTS，以免不太符合本库场景或已废弃等原因造成你的贡献浪费，那就太可惜了！

## 2.1.2 示例蜘蛛的演练

为了向您展示 ayugespidertools 带来了什么，我们将带您通过一个 Scrapy Spider 示例，使用最简单的方式来运行蜘蛛。

先创建项目：

```
# eg: 本示例使用的 project_name 为 DemoSpider  
  
ayuge startproject <project_name>
```

创建爬虫脚本：

```
进入项目根目录  
cd <project_name>  
  
生成脚本  
ayuge genspider <spider_name> <example.com>
```

下面是从网站 <https://blog.csdn.net/phoenix/web/blog/hot-rank?page=0&pageSize=25&type=> 抓取热榜信息的蜘蛛代码：

```
import json  
  
from ayugespidertools.common.utils import ToolsForAyu  
from ayugespidertools.items import DataItem, MysqlDataItem  
from ayugespidertools.spiders import AyuSpider  
from scrapy.http import Request  
from scrapy.http.response.text import TextResponse  
  
from DemoSpider.items import TableEnum  
from DemoSpider.settings import logger  
  
"""  
#####  
↪#####  
# collection_website: CSDN - 专业开发者社区  
# collection_content: 热榜文章排名 Demo 采集示例 - 存入 Mysql (配置根据本地 settings 的_  
↪LOCAL_MYSQL_CONFIG 中取值)  
# create_time: 2022-07-30  
# explain:  
# demand_code_prefix = ""  
#####  
↪#####  
"""
```

(continues on next page)

(continued from previous page)

```

class DemoOneSpider(AyuSpider):
    name = "demo_one"
    allowed_domains = ["blog.csdn.net"]
    start_urls = ["https://blog.csdn.net/"]

    # 数据库表的枚举信息, 当前项目所依赖的表信息, 一般用于存储数据时使用
    custom_table_enum = TableEnum
    # 初始化配置的类型, 初始化设置 (不需要直接不用配置)
    settings_type = 'debug'
    custom_settings = {
        # scrapy 日志等级
        'LOG_LEVEL': 'DEBUG',
        # 是否开启记录项目相关运行统计信息。不配置默认为 False
        "RECORD_LOG_TO_MYSQL": True,
        # 设置 ayugespidertools 库的日志输出为 loguru, 可自行配置 logger 规则来管理项目日志。若不配置此项, 库日志只会在控制台上打印
        "LOGURU_CONFIG": logger,
        # Mysql 数据表的前缀名称, 用于标记属于哪个项目, 也可以不用配置
        "MYSQL_TABLE_PREFIX": "demo1_",
        "ITEM_PIPELINES": {
            # 激活此项则数据会存储至 Mysql
            "ayugespidertools.pipelines.AyuFtyMysqlPipeline": 300,
        },
        "DOWNLOADER_MIDDLEWARES": {
            # 随机请求头
            "ayugespidertools.middlewares.RandomRequestUaMiddleware": 400,
        },
    }

    # 打开 mysql 引擎开关, 用于数据入库前更新逻辑判断
    mysql_engine_enabled = True

    def start_requests(self):
        """
        获取项目热榜的列表数据
        """
        yield Request(
            url="https://blog.csdn.net/phoenix/web/blog/hot-rank?page=0&pageSize=25&
↪type=",
            callback=self.parse_first,
            headers={

```

(continues on next page)

(continued from previous page)

```
        "referer": "https://blog.csdn.net/rank/list",
    },
    dont_filter=True,
)

def parse_first(self, response: TextResponse):
    data_list = json.loads(response.text) ["data"]
    for curr_data in data_list:
        # 这里的所有解析方式可选择自己习惯的其它任意库, xpath, json 或正则等等。
        article_detail_url = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="articleDetailUrl"
        )

        article_title = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="articleTitle"
        )

        comment_count = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="commentCount"
        )

        favor_count = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="favorCount"
        )

        nick_name = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="nickName"
        )

        # 数据存储方式 1, 需要添加注释时的写法
        ArticleInfoItem = MysqlDataItem(
            # 这里也可以写为 article_detail_url = DataItem(article_detail_url), 但没有
            # 功能了, 那不如使用下面的数据存储方式 2
            article_detail_url=DataItem(article_detail_url, "文章详情链接"),
            article_title=DataItem(article_title, "文章标题"),
            comment_count=DataItem(comment_count, "文章评论数量"),
            favor_count=DataItem(favor_count, "文章赞成数量"),
            nick_name=DataItem(nick_name, "文章作者昵称"),
            _table=TableEnum.article_list_table.value["value"],
        )

        # 数据存储方式 2, 若不需要注释, 也可以这样写, 但不要两种风格混用
```

注释

(continues on next page)

(continued from previous page)

```

"""
ArticleInfoItem = MysqlDataItem(
    article_detail_url=article_detail_url,
    article_title=article_title,
    comment_count=comment_count,
    favor_count=favor_count,
    nick_name=nick_name,
    _table=TableEnum.article_list_table.value["value"],
)
"""

self.slog.info(f"ArticleInfoItem: {ArticleInfoItem}")
# yield ArticleInfoItem

# 数据入库逻辑 -> 测试 mysql_engine 的去重功能, 你可以自行实现。mysql_engine 也已经
给你了

save_table = f'{self.custom_settings.get("MYSQL_TABLE_PREFIX", "")}'
->{TableEnum.article_list_table.value["value"]}
sql = f"""select `id` from `{save_table}` where `article_detail_url` = "
->{article_detail_url}" limit 1"""
yield ToolsForAyu.filter_data_before_yield(
    sql=sql,
    mysql_engine=self.mysql_engine,
    item=ArticleInfoItem,
)

```

## 刚刚发生了什么？

刚刚使用 ayugespidertools 创建了项目，并生成了具体的爬虫脚本示例。其爬虫脚本中的各种依赖（比如项目目录结构，配置信息等）在创建项目后就正常产生了，一般所需的配置信息（比如 MySQL, MongoDB, OSS 等）在项目的 VIT 目录下 .conf 文件中修改，不需要配置的不用理会它即可。

只要配置好 .conf 信息，就可以跑通以上示例。如果修改为新的项目，只需要修改上面示例中的 spider 解析规则即可。

### 2.1.3 还有什么？

本库依赖 Scrapy，你可以使用 Scrapy 命令来管理你的项目，体会 Scrapy 的强大和方便。

ayugespidertools 根据 scrapy 的模板功能方便的创建示例脚本，比如：

```

# 查看支持的脚本模板示例
ayuge genspider -l

```

(continues on next page)

(continued from previous page)

```
<output>
Available templates:
  async
  basic
  crawl
  csvfeed
  xmlfeed

# 使用具体的示例命令
ayuge genspider -t <Available_templates> <spider_name> <example.com>

eg: ayuge gendpier -t async demom_async baidu.com
```

### 2.1.4 下一步是什么？

接下来的步骤是 安装 AyugeSpiderTools，按照 Scrapy 的教程学习如何使用 Scrapy 并加入 Scrapy 社区。谢谢你的关注！

## 2.2 安装指南

### 2.2.1 支持的 Python 版本

AyugeSpiderTools 需要 Python 3.8.1+。

### 2.2.2 安装 AyugeSpiderTools

为了向您展示 ayugespidertools 带来了什么，我们将带您通过一个 Scrapy Spider 示例，使用最简单的方式来运行蜘蛛。

可以使用以下命令从 PyPI 安装 ayugespidertools 及其依赖项：

```
pip install ayugespidertools
```

我们强烈建议您将 ayugespidertools 安装在专用的 virtualenv 中，以避免与您的系统包发生冲突。

## 值得知道的事情

ayugespidertools 是依赖于 Scrapy 开发的，对其在爬虫开发中遇到的常用操作进行扩展。

## 使用虚拟环境（推荐）

建议在所有平台上的虚拟环境中安装此库。

有关如何创建虚拟环境的信息，请参阅[虚拟环境和包](#)。

## 2.3 AyugeSpiderTools 教程

在本教程中，我们假设您的系统上已经安装了 ayugespidertools。

我们要抓取 [blog.csdn.net](http://blog.csdn.net)，这是一个知识问答社区的网站。

本教程将引导您完成这些任务：

- 创建一个新的 Scrapy 项目
- 编写爬虫来抓取站点并提取数据
- 使用命令行导出抓取的数据
- 更改蜘蛛以递归地跟踪链接
- 使用蜘蛛参数

### 2.3.1 创建项目

在开始抓取之前，您必须设置一个新的 ayugespidertools 项目。输入您要存储代码的目录并运行：

```
ayuge startproject DemoSpider
```

这将创建一个 DemoSpider 包含以下内容的目录：

```
DemoSpider/
|-- DemoSpider                # project's Python module, you'll import your_
↪code from here
|   |-- __init__.py
|   |-- items.py              # project items definition file, 数据库表枚举信息示例也迁移至此
|   |-- logs                  # 日志管理文件夹, 可以自定义规则
|   |   |-- DemoSpider.log    # scrapy 输出日志, 文件名称为项目名
|   |   |-- error.log         # loguru 日志 error 规则输出文件
|   |-- middlewares.py        # project middlewares definition file
```

(continues on next page)

(continued from previous page)

```

| |-- pipelines.py           # project pipelines definition file
| |-- run.py                # scrapy 运行文件
| |-- run.sh                # 项目运行 shell, 运行以上的 run.py, win 平台不会生成此
文件
| |-- settings.py          # project settings definition file
| |-- spiders               # a directory where you'll later put your
↳spiders
| | |-- __init__.py
| | `-- VIT
|     |-- .conf             # 配置文件, 用于修改 Mysql, MongoDB 等配置
|-- pyproject.toml         # 项目配置
|-- README.md              # 说明文档
|-- requirements.txt       # 依赖文件
`-- scrapy.cfg             # deploy configuration file

```

## 2.3.2 我们的第一个 Spider

这是我们第一个 Spider 的代码。demo\_one.py 将其保存在项目目录下命名的文件 DemoSpider/spiders 中:

```

from ayugespidertools.common.utils import ToolsForAyu
from ayugespidertools.items import DataItem, MongoDataItem, MysqlDataItem
from ayugespidertools.spiders import AyuSpider
from scrapy.http import Request

from DemoSpider.items import TableEnum

"""
#####
↳#####
# collection_website: CSDN - 专业开发者社区
# collection_content: 热榜文章排名 Demo 采集示例 - 同时存入 Mysql 和 MongoDB 的场景
# create_time: 2022-08-22
# explain: 根据本项目中的 demo_one 脚本修改而得
# demand_code_prefix = ""
#####
↳#####
"""

class DemoEightSpider(AyuSpider):
    name = "demo_eight"

```

(continues on next page)



(continued from previous page)

```

allowed_domains = ["blog.csdn.net"]
start_urls = ["https://blog.csdn.net/"]

# 数据库表的枚举信息
custom_table_enum = TableEnum
# 初始化配置的类型
settings_type = "debug"
custom_settings = {
    # 是否开启记录项目相关运行统计信息
    "RECORD_LOG_TO_MYSQL": False,
    # 数据表的前缀名称, 用于标记属于哪个项目, 也可以不用添加
    "MYSQL_TABLE_PREFIX": "demo8_",
    # 数据表的前缀名称, 用于标记属于哪个项目 (也可不配置此参数, 按需配置)
    "MONGODB_COLLECTION_PREFIX": "demo8_",
    "ITEM_PIPELINES": {
        # 激活此项则数据会存储至 Mysql
        "ayugespidertools.pipelines.AyuFtyMysqlPipeline": 300,
        # 激活此项则数据会存储至 MongoDB
        "ayugespidertools.pipelines.AyuFtyMongoPipeline": 301,
    },
    "DOWNLOADER_MIDDLEWARES": {
        # 随机请求头
        "ayugespidertools.middlewares.RandomRequestUaMiddleware": 400,
    },
}

# 打开 mysql 引擎开关, 用于数据入库前更新逻辑判断
mysql_engine_enabled = True

def start_requests(self):
    """
    get 请求首页, 获取项目列表数据
    """
    yield Request(
        url="https://blog.csdn.net/phoenix/web/blog/hot-rank?page=0&pageSize=25&
↪type=",
        callback=self.parse_first,
        headers={
            "referer": "https://blog.csdn.net/rank/list",
            "user-agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/
↪537.36 (KHTML, like Gecko) Chrome/103.0.0.0 Safari/537.36",
        },
        dont_filter=True,

```

(continues on next page)

(continued from previous page)

```
)

def parse_first(self, response):
    data_list = ToolsForAyu.extract_with_json(
        json_data=response.json(), query="data"
    )
    for curr_data in data_list:
        article_detail_url = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="articleDetailUrl"
        )

        article_title = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="articleTitle"
        )

        comment_count = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="commentCount"
        )

        favor_count = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="favorCount"
        )

        nick_name = ToolsForAyu.extract_with_json(
            json_data=curr_data, query="nickName"
        )

        ArticleMysqlItem = MysqlDataItem(
            article_detail_url=DataItem(article_detail_url, "文章详情链接"),
            article_title=DataItem(article_title, "文章标题"),
            comment_count=DataItem(comment_count, "文章评论数量"),
            favor_count=DataItem(favor_count, "文章赞成数量"),
            nick_name=DataItem(nick_name, "文章作者昵称"),
            _table=TableEnum.article_list_table.value["value"],
        )
        yield ArticleMysqlItem

        ArticleMongoItem = MongoDataItem(
            article_detail_url=article_detail_url,
            article_title=article_title,
            comment_count=comment_count,
            favor_count=favor_count,
            nick_name=nick_name,
```

(continues on next page)

(continued from previous page)

```

        _table=TableEnum.article_list_table.value["value"],
        # 这里表示以 article_detail_url 为去重规则, 若存在则更新, 不存在则新增
        _mongo_update_rule={"article_detail_url": article_detail_url},
    )
    yield ArticleMongoItem

```

如您所见, 我们的 Spider 子类化 `AyugeSpider.AyuSpider` 并定义了一些属性和方法:

- `name`: 标识蜘蛛。在一个项目中必须是唯一的, 即不能为不同的 Spiders 设置相同的名字。
- `start_requests()`: 必须返回一个可迭代的请求 (你可以返回一个请求列表或编写一个生成器函数), 蜘蛛将从中开始爬行。后续请求将从这些初始请求中依次生成。
- `parse_first()`: 将被调用以处理为每个请求下载的响应的方法。`response` 参数是 `TextResponse` 的一个实例, 它保存页面内容, 并有进一步的有用方法来处理它。该 `parse_first()` 方法通常解析响应, 将抓取的数据提取为字典, 并找到要遵循的新 URL 并从中创建新请求 (`Request`)。

另外, 一些其它注意事项:

- 示例中的一些配置和一些功能并不是每个项目中都必须要编写和配置的, 只是用于展示一些功能
- 据上条可知, 可以写出很简洁的代码, 删除你认为的无关配置和方法并将其配置成你自己的模板就更容易适配更多人的使用场景。

## 如何运行我们的蜘蛛

要让我们的蜘蛛工作, 请转到项目的顶级目录并运行:

```

# 本身就是 scrapy 的项目, 所以可以使用 scrapy 可以执行的任何形式即可
scrapy crawl demo_one

# 或者执行项目根目录下的 run.py (需要编辑自己需要执行的脚本)
python run.py

# 或者执行项目根目录下的 run.sh, 其实它也是通过调用 run.py 来执行的。只不过 shell 文件中包含了虚拟环境的 activate 了而已
sh run.sh

```

## 2.4 例子

本教程将引导您完成这些任务：

- 快速熟悉 ayugespidertools 库的使用方法和支持场景
- 编写爬虫来抓取站点并提取数据

### 2.4.1 1. 快速开始

你可以使用以下两种方式来快速开始

#### 1.1. 方式一：ayugespidertools

通过跑通本库 Github 中的 GIF 示例

具体请点击跳转至 [AyugeSpiderTools](#) 查看

#### 1.2. 方式二：DemoSpider

通过另一个的演示项目 DemoSpider 来选择复现某些场景

最好的学习方法是通过 Github 上的 DemoSpider 示例，您可以使用它快速复现某些场景下的功能。

本库 ayugespidertools 的 [github README.md](#) 中所有功能，都可以在 DemoSpider 中找到示例。

DemoSpider 项目位于：<https://github.com/shengchenyang/DemoSpider>，您可以在项目的自述文件中找到有关它的更多信息。

### 2.4.2 2. 应用场景介绍

根据 DemoSpider 中的各个 spider 对一些应用场景进行简要的补充介绍，总体的介绍为：

```
# 采集数据存入 `Mysql` 的场景：
+ 1) .demo_one: 配置根据本地 `settings` 的 `LOCAL_MYSQL_CONFIG` 中取值
+ 3) .demo_three: 配置根据 `consul` 的应用管理中心中取值
+ 5) .demo_five: 异步存入 `Mysql` 的场景

# 采集数据存入 `MongoDB` 的场景：
+ 2) .demo_two: 采集数据存入 `MongoDB` 的场景（配置根据本地 `settings` 的 `LOCAL_MONGODB_
↪CONFIG` 中取值）
+ 4) .demo_four: 采集数据存入 `MongoDB` 的场景（配置根据 `consul` 的应用管理中心中取值）
+ 6) .demo_six: 异步存入 `MongoDB` 的场景
```

(continues on next page)

(continued from previous page)

```

# 将 `Scrapy` 的 `Request`, `FormRequest` 替换为其它工具实现的场景
- 以上为使用 scrapy Request 的场景
+ 7).demo_seven: scrapy Request 替换为 requests 请求的场景 (一般情况下不推荐使用, 同步库
+ 会拖慢 scrapy 速度, 可用于测试场景)

+ 8).demo_eight: 同时存入 Mysql 和 MongoDB 的场景

+ 9).demo_aiohttp_example: scrapy Request 替换为 aiohttp 请求的场景, 提供了各种请求场景示例
(GET,POST)
+ 10).demo_aiohttp_test: scrapy aiohttp 在具体项目中的使用方法示例

+ 11).demo_proxy_one: 快代理动态隧道代理示例
+ 12).demo_proxy_two: 测试快代理独享代理

+13).demo_AyuTurboMysqlPipeline: mysql 同步连接池的示例
+14).demo_crawl: 支持 scrapy CrawlSpider 的示例

# 本库中给出支持 Item Loaders 特性的示例 (文档地址: https://ayugespidertools.readthedocs.io/
->en/latest/topics/loaders.html)
+15).demo_item_loader: 本库中使用 Item Loaders 的示例
-16).demo_item_loader_two: 展示本库使用 itemLoader 特性的示例, 此示例已删除, 可查看上个 demo_
->item_loader 中的示例, 目标已经可以很方便的使用 Item Loaders 功能了

+17).demo_mongo_async: asyncio 版本存储 mongoDB 的 pipelines 示例

```

基本查看查看以上 spider 即可了解使用方法, 但有些示例还是不够详细, 对以上内容进行补充。

- 以上场景有需要 consul 上的相关配置的示例, 是根据 .conf 中的配置并按照小写风格书写, 以下为 json 格式配置的示例:

```

{
  "mysql": {
    "host": "****",
    "port": 3306,
    "user": "****",
    "password": "****",
    "database": "****",
    "charset": "utf8mb4"
  },
  "mongodb": {
    "host": "****",
    "port": 27017,
    "user": "****",
    "password": "****",

```

(continues on next page)

```
    "database": "***",
    "authsource": "***"
  },
  "consul": {
    "token": "",
    "url": "http://host:port/v1/kv/...?dc=dc1",
    "format": "json"
  },
  "kdl_dynamic_proxy": {
    "proxy": "o668.kdltps.com:15818",
    "username": "***",
    "password": "***"
  },
  "kdl_exclusive_proxy": {
    "proxy": "http://kps.kdlapi.com/api/getkps?orderid=***&num=100&format=json
↪",
    "username": "***",
    "password": "***",
    "index": 1
  },
  "ali_oss": {
    "accesskeyid": "LTA***",
    "accesskeysecret": "***",
    "endpoint": "https://oss-cn-***.aliyuncs.com",
    "bucket": "***",
    "doc": "***"
  }
}
```

## *AyugeSpiderTools* 一目了然

了解什么是 *AyugeSpiderTools* 以及它如何为您提供帮助。

## 安装指南

在您的设备上安装 *AyugeSpiderTools*。

## *AyugeSpiderTools* 教程

编写您的第一个 *AyugeSpiderTools* 项目。

## 例子

通过一个简单示例来了解一些信息。

## 基本概念

### 3.1 命令行工具

AyugeSpiderTools 是直接使用 scrapy 命令行工具来控制的，这里简称为 AyugeSpiderTools 工具，以区别于我们简称为“命令”或“Scrapy 命令”的子命令。

AyugeSpiderTools 工具只提供了常用的几个命令，用于多种用途，每个命令都接受一组不同的参数和选项。但是，其它缺失的命令你可以直接使用 Scrapy 的即可。

#### 3.1.1 配置设置

未改变，也是在在标准位置的 ini 样式文件中查找配置参数 scrapy.cfg:

这些文件中的设置按列出的优先顺序合并：用户定义的值比系统范围的默认值具有更高的优先级，并且项目范围的设置将在定义时覆盖所有其他设置。

#### 3.1.2 AyugeSpiderTools 项目的默认结构

在深入研究命令行工具及其子命令之前，让我们先了解一下项目的目录结构。

虽然可以修改，但是所有的项目默认都有相同的文件结构，类似这样：

```
|-- DemoSpider
|   |-- __init__.py
|   |-- items.py
|   |-- logs
|   |   |-- DemoSpider.log
|   |   |-- error.log
|   |-- middlewares.py
|   |-- pipelines.py
|   |-- run.py
|   |-- run.sh
|   |-- settings.py
```

(continues on next page)

(continued from previous page)

```
| |-- spiders
| | |-- __init__.py
| | `-- spider1.py
| `-- VIT
|     |-- .conf
|-- pyproject.toml
|-- README.md
|-- requirements.txt
`-- scrapy.cfg
```

文件所在的目录 `scrapy.cfg` 称为项目根目录。该文件包含定义项目设置的 python 模块的名称。这是一个例子：

```
[settings]
default = DemoSpider.settings
```

### 3.1.3 使用 Ayugespidertools 工具

您可以先运行不带参数的 AyugeSpiderTools 工具，它会打印一些使用帮助和可用的命令：

命令如下：

```
ayuge -h
```

输出示例如下：

```
AyugeSpiderTools 3.0.1 - no active project

Usage:
  ayuge <command> [options] [args]

Available commands:
  genspider      Generate new spider using pre-defined templates
  startproject  Create new project
  version       Print AyugeSpiderTools version

  [ more ]      More commands available when run from project directory

Use "ayuge <command> -h" to see more info about a command
```



## 创建项目

您通常使用该工具做的第一件事是创建您的项目：

```
ayuge startproject myproject [project_dir]
```

这将在该目录下创建一个 Scrapy 项目 `project_dir`。如果 `project_dir` 未指定，则项目目录将与 `myproject` 相同。

接下来，进入新项目目录：

```
cd project_dir
```

您已准备好使用该命令从那里管理和控制您的项目。

## 控制项目

您可以使用项目内部的工具来控制和管理它们。

例如，要创建一个新的蜘蛛：

```
ayuge genspider mydomain mydomain.com
```

### 3.1.4 启动项目

- 句法: `ayuge startproject <project_name> [project_dir]`
- 需要项目: 无

在 `project_dir` 目录下创建一个名为 `project_name` 的新项目。如果未指定项目目录，则项目目录将与项目名称相同。

使用示例：

```
ayuge startproject myproject
```

### 3.1.5 可用的工具命令

本节包含可用内置命令的列表以及说明和一些用法示例。请记住，您始终可以通过运行以下命令获取有关每个命令的更多信息：

```
ayuge <command> -h
```

您可以使用以下命令查看所有可用命令：

```
ayuge -h
```

## 启动项目

- 句法: `ayuge startproject <project_name> [project_dir]`
- 需要项目: 无

在 `project_dir` 目录下创建一个名为 `project name` 的新项目。如果未指定项目目录, 则项目目录将与项目名称相同。

使用示例:

```
ayuge startproject myproject
```

## genspider

- 句法: `ayuge genspider [-t template] <name> <domain or URL>`
- 需要项目: 无

使用示例:

```
$ ayuge genspider -l
Available templates:
  async
  basic
  crawl
  csvfeed
  xmlfeed

$ ayuge genspider example example.com
Created spider 'example' using template 'basic'

$ ayuge genspider -t crawl scrapyorg scrapy.org
Created spider 'scrapyorg' using template 'crawl'
```

## 3.2 Item

演示在使用本库场景下 Item 的使用方法。

本教程将引导您完成这些任务：

- 演示本库推荐的 Item 适配方式
- 补充适配 `add_value`, `add_xpath`, `add_css` 等方法的示例

### 3.2.1 实现原理

以下为本库中推荐的 `mysql` 和 `MongoDB` 存储时的主要 Item 示例：

其实,本库就是推荐把所有字段统一存入 `alldata` 字段中,其它字段用于场景补充,比如:`table` 字段用于说明要存储的表名/集合名,`item_mode` 字段用于说明存储的方式,`mongo_update_rule` 字段是 `item_mode` 为 `MongoDB` 存储场景时的去重条件(可不设置此字段)。

本库,

```
def parse(self, response):
    # 存储到 Mysql 场景时需要的 Item 构建示例
    ArticleMysqlItem = MysqlDataItem(
        article_detail_url=DataItem(article_detail_url, "文章详情链接"),
        article_title=DataItem(article_title, "文章标题"),
        comment_count=DataItem(comment_count, "文章评论数量"),
        favor_count=DataItem(favor_count, "文章赞成数量"),
        nick_name=DataItem(nick_name, "文章作者昵称"),
        _table=TableEnum.article_list_table.value["value"],
    )
    # 存储到 MongoDB 场景时需要的 Item 构建示例
    ArticleMongoItem = MongoDataItem(
        article_detail_url=article_detail_url,
        article_title=article_title,
        comment_count=comment_count,
        favor_count=favor_count,
        nick_name=nick_name,
        _table=TableEnum.article_list_table.value["value"],
        # 这里表示以 article_detail_url 为去重规则,若存在则更新,不存在则新增
        _mongo_update_rule={"article_detail_url": article_detail_url},
    )
```

以上可知,目前可直接将需要的参数在对应 Item 中直接按 `key=value` 赋值即可, `key` 即为存储至库中字段, `value` 为存储内容。

当然,目前也支持动态赋值,但是我不推荐使用,直接按照上个方式即可:

```
def parse(self, response):
    mdi = MysqlDataItem(_table="table0")
    mdi.add_field("add_field1", "value1")
    mdi.add_field("add_field2", DataItem(key_value="value2"))
    mdi.add_field("add_field3", DataItem(key_value="value3", notes="add_field3 值"))
    # _table 修改可通过以下方式, 同样不推荐使用
    mdi._table = "table1"
```

另外, 本库的 item 提供类型转换, 以方便后续的各种使用场景:

```
# 将本库 Item 转为 dict 的方法
item_dict = mdi.asdict()
# 将本库 Item 转为 scrapy Item 的方法
item = mdi.asitem()
```

### 3.2.2 使用示例

只需要在 yield item 时, 按需提前导入 MysqlDataItem, MongoDataItem, 将所有的存储字段和场景补充字段全部添加完整即可。

以本库模板中的 basic.tpl 为例:

```
#!/usr/bin/env python3
# -*- coding:utf-8 -*-
import pandas
from DemoSpider.settings import logger
from scrapy.http.response.text import TextResponse

from ayugespidertools.common.utils import ToolsForAyu
from ayugespidertools.items import DataItem, MongoDataItem, MysqlDataItem
from ayugespidertools.spiders import AyuSpider
from scrapy.http import Request

from DemoSpider.items import TableEnum

"""
#####
↪#####
# collection_website: csdn.net - 采集的目标站点介绍
# collection_content: 采集内容介绍
# create_time: xxxx-xx-xx
# explain:
# demand_code_prefix = ''
```

(continues on next page)

(continued from previous page)

```
#####
->#####
"""

class DemoOneSpider(AyuSpider):
    name = 'demo_one'
    allowed_domains = ['csdn.net']
    start_urls = ['https://www.csdn.net/']

    # 数据库表的枚举信息
    custom_table_enum = TableEnum
    # 初始化配置的类型
    settings_type = 'debug'
    custom_settings = {
        # scrapy 日志等级配置
        'LOG_LEVEL': 'DEBUG',
        # 以 loguru 来管理日志, 本库会在 settings 中生成规则示例, 可自行修改。也可不配置
        'LOGURU_CONFIG': logger,
        # Mysql 数据表的前缀名称, 用于标记属于哪个项目, 可不配置
        'MYSQL_TABLE_PREFIX': "demo_basic_",
        # MongoDB 集合的前缀名称, 用于标记属于哪个项目, 可不配置
        'MONGODB_COLLECTION_PREFIX': "demo_basic_",
        'ITEM_PIPELINES': {
            # 激活此项则数据会存储至 Mysql
            'ayugespidertools.pipelines.AyuFtyMysqlPipeline': 300,
            # 激活此项则数据会存储至 MongoDB
            'ayugespidertools.pipelines.AyuFtyMongoPipeline': 301,
        },
        'DOWNLOADER_MIDDLEWARES': {
            # 随机请求头
            'ayugespidertools.middlewares.RandomRequestUaMiddleware': 400,
        },
    }

    # 打开 mysql 引擎开关, 用于数据入库前更新逻辑判断
    mysql_engine_enabled = True

    def start_requests(self):
        """
        get 请求首页, 获取项目列表数据
        """
        yield Request(
```

(continues on next page)

(continued from previous page)

```

url="https://blog.csdn.net/phoenix/web/blog/hot-rank?page=0&pageSize=25&
↪type=",
callback=self.parse_first,
headers={
    'referer': 'https://blog.csdn.net/rank/list',
},
cb_kwargs={
    "curr_site": "csdn",
},
dont_filter=True
)

def parse_first(self, response: TextResponse, curr_site: str):
    # 日志使用: scrapy 的 self.logger 或本库的 self.slog 或直接使用全局的 logger handle_
    ↪也行 (根据场景自行选择)
    self.slog.info(f"当前采集的站点为: {curr_site}")

    # 你可以自定义解析规则, 使用 lxml 还是 response.css response.xpath 等等都可以。
    data_list = ToolsForAyu.extract_with_json(json_data=response.json(), query=
    ↪"data")
    for curr_data in data_list:
        article_detail_url = ToolsForAyu.extract_with_json(
            json_data=curr_data,
            query="articleDetailUrl")

        article_title = ToolsForAyu.extract_with_json(
            json_data=curr_data,
            query="articleTitle")

        comment_count = ToolsForAyu.extract_with_json(
            json_data=curr_data,
            query="commentCount")

        favor_count = ToolsForAyu.extract_with_json(
            json_data=curr_data,
            query="favorCount")

        nick_name = ToolsForAyu.extract_with_json(
            json_data=curr_data,
            query="nickName")

        ArticleMysqlItem = MysqlDataItem(
            article_detail_url=DataItem(article_detail_url, "文章详情链接"),

```

(continues on next page)

(continued from previous page)

```

        article_title=DataItem(article_title, "文章标题"),
        comment_count=DataItem(comment_count, "文章评论数量"),
        favor_count=DataItem(favor_count, "文章赞成数量"),
        nick_name=DataItem(nick_name, "文章作者昵称"),
        _table=TableEnum.article_list_table.value["value"],
    )
    yield ArticleMysqlItem

    ArticleMongoItem = MongoDataItem(
        article_detail_url=article_detail_url,
        article_title=article_title,
        comment_count=comment_count,
        favor_count=favor_count,
        nick_name=nick_name,
        _table=TableEnum.article_list_table.value["value"],
        # 这里表示以 article_detail_url 为去重规则, 若存在则更新, 不存在则新增
        _mongo_update_rule={"article_detail_url": article_detail_url},
    )
    yield ArticleMongoItem

```

由上可知，本库中的 Item 使用方法还是很方便的。

#### 对以上 Item 相关信息解释：

- 先导入所需 Item
  - mysql 场景导入 MysqlDataItem
  - mongo 场景导入 MongoDataItem
- 构建对应场景的 Item
  - MysqlDataItem 构建 Mysql 存储场景
  - MongoDataItem 构建 MongoDB 存储场景
- 最后 yield 对应 item 即可

### 3.2.3 yield item

这里解释下 item 的格式问题，虽说也是支持直接 yield dict，scrapy 的 item 格式 (即 ScrapyClassicItem)，还有就是本库推荐的 MysqlDataItem 和 MongoDataItem 的形式。

这里介绍下 item 字段及其注释，以上所有 item 都有参数提示：

注，对以上表格中内容进行扩充解释：

- 自定义字段使用示例请在 [readthedocs](#) 中查看。

### 3.2.4 自定义 Item 字段和实现 Item Loaders

本库支持 scrapy Item 的格式，或者 dict 格式，DemoSpider 中 demo\_one 已有示例。

由于本库推荐将所有存储字段统一存储至 alldata 中来维护的，而且本库 Item 已使用了 scrapy 推荐的 dataclass 实现，自然是不推荐自定义 Item 字段的，但是也是可以实现的。具体请在下一章浏览。

## 3.3 Item Loaders

Item Loaders 提供了一种方便的机制来填充已抓取的 ITEM。尽管项目可以直接填充，项目加载器提供了一个更方便的 API 来从抓取过程中填充它们，通过自动化一些常见的任务，比如在分配它之前解析原始提取的数据。

换句话说，ITEM 提供了抓取数据的容器，而项目加载器提供了**填充**该容器的机制。

Item Loaders 旨在提供一种灵活、高效和简单的机制来扩展和覆盖不同的字段解析规则，无论是通过蜘蛛，还是通过源格式 (HTML、XML 等)，而不会成为维护的噩梦。

具体请查看 scrapy 中对应的 [Item Loaders](#) 的文档。

由文档可知，如果使用 Item Loaders 需要先声明 Item 子类，并固定 Field 字段。即以下内容示例：

```
import scrapy

class Product(scrapy.Item):
    book_name = scrapy.Field()
    book_href = scrapy.Field()
    book_intro = scrapy.Field()
```

但是本库不固定 Item field 的内容，这样丧失了解放双手的目的。

虽然，scrapy 也可以通过使用如下方法来新增字段，但总归 scrapy 是不推荐这样的写法且不太方便：

```
Product.fields["add_field1"] = scrapy.Field()
```

那本库如何实现使用项目加载器填充项目的效果呢，本库是通过 Item 的 asitem 方法实现。具体使用方法请看文章后半段。



### 3.3.1 使用项目加载器填充项目

这是 Spider 中典型的 Item Loader 用法:

```
from scrapy.loader import ItemLoader
from myproject.items import Product
from ayugespidertools.items import MongoDataItem, MysqlDataItem

# 这是 scrapy 中的实现示例:
def parse(self, response):
    l = ItemLoader(item=Product(), response=response)
    l.add_xpath('name', '//div[@class="product_name"]')
    l.add_xpath('name', '//div[@class="product_title"]')
    l.add_xpath('price', '//p[@id="price"]')
    l.add_css('stock', 'p#stock')
    l.add_value('last_updated', 'today') # you can also use literal values
    yield l.load_item()

# 这是本库中的实现示例:
def parse(self, response):
    # 先定义所需字段
    my_product = MysqlDataItem(
        book_name=None,
        book_href=None,
        book_intro=None,
        _table=TableEnum.article_list_table.value["value"],
    )
    # 然后可使用 asitem 的方法使用常规的 ItemLoader 功能
    mine_item = ItemLoader(item=my_product.asitem(), selector=None)
    mine_item.default_output_processor = TakeFirst()
    mine_item.add_value("book_name", book_name)
    mine_item.add_xpath("book_href", '//div[@class="product_title"]')
    mine_item.add_css("book_intro", 'p#stock')
    item = mine_item.load_item()
    yield item
```

### 3.3.2 使用数据类项

那本库的方式在使用 `ItemLoader` 时有没有缺点呢？

是的，有缺点，由于本库虽然支持动态添加 `Item` 字段，但是其实不太好实现 `dataclass items` 的字段类型约束和参数 `default` 的相关设置。本库是不推荐固定 `Item` 字段（比如 `ayugespidertools v3.0.0` 之前的版本中，会把数据都存入 `alldata` 的固定字段中），也不推荐不同 `spider` 就需要定义其对应的不同 `Item class` 的。其实各有优缺点，只是本库选择了牺牲此部分。

### 3.3.3 输入和输出处理器

那么，在本库中的使用方法如下：

`Item Loader` 包含一个输入处理器和一个用于每个 `item` 字段的输出处理器。输入处理器在收到数据后立即处理提取的数据（通过 `add_xpath()`, `add_css()` 或 `add_value()` 方法），输入处理器的结果被收集并保存在 `ItemLoader` 中。收集完所有数据后，`ItemLoader.load_item()` 调用该方法填充并获取填充的项对象。那是使用先前收集的数据（并使用输入处理器处理）调用输出处理器的时候。输出处理器的结果是分配给项目的最终值。

让我们看一个示例来说明如何为特定字段调用输入和输出处理器（这同样适用于任何其他字段）：

```
l = ItemLoader(my_product.asitem(), some_selector)
l.default_output_processor = TakeFirst()
l.add_xpath("name", xpath1) # (1)
l.add_xpath("name", xpath2) # (2)
l.add_css("name", css) # (3)
l.add_value("name", "test") # (4)
return l.load_item() # (5)
```

然后就可以使用 [使用项目加载器填充项目](#使用项目加载器填充项目) 中的代码了

本库主推便捷，不太推荐使用以上代码自定义增加 `Item` 字段来适配 `Item Loaders` 的特性，除非某些场景下使用 `Item Loaders` 能够极大方便开发时，才推荐使用下。

#### 命令行工具

了解用于管理 `Scrapy` 项目的命令行工具。

#### *Item*

定义要采集的数据。

#### *Item Loaders*

用提取的数据填充你的 `item`。

## 4.1 Logging

Scrapy 用 `logging` 来记录日志，日志记录开箱即用，并且可以在某种程度上使用日志设置中列出的 Scrapy 设置进行配置。

本文不再介绍其详细配置及用法，请移步其官网文档中查看。AyugeSpiderTools 库会在 `settings` 中默认设置一个日志存储配置，默认放在项目的 `logs` 文件夹下，其名称为项目名称，如下所示：

```
# 日志管理
LOG_FILE = f"{LOG_DIR}/DemoSpider.log"

# 配置中 DemoSpider 是与 ayugespidertools startproject <project_name> 中的项目名称对应的
```

### 4.1.1 slog 日志

`ayugespidertools` 添加了 `loguru` 库来管理日志，可以很方便的查看不同日志等级的信息。可以在 `spider` 脚本中使用 `spider.slog` 或 `self.slog` 即可记录日志。同样，本库会在 `settings` 中也默认设置一个 `loguru` 的日志存储配置，也放在项目的 `logs` 文件夹下，如下所示：

```
# 本库只会持久化记录 error 级别的日志，但在调试时也可以方便地查看（包括其它等级的）控制台日志
logger.add(
    env.str("LOG_ERROR_FILE", f"{LOG_DIR}/error.log"),
    level="ERROR",
    rotation="1 week",
    retention="7 days",
    enqueue=True,
)
```

## 4.1.2 日志级别

ayugespidertools 中的 loguru 日志等级与 Python 的内置日志记录定义一致，大致分为 5 个不同的级别来指示给定日志消息的严重性。以下是标准的，按降序排列：

1. `slog.CRITICAL`- 对于严重错误（最高严重性）
2. `slog.ERROR`- 对于常规错误
3. `slog.WARNING`- 警告信息
4. `slog.INFO`- 用于信息性消息
5. `slog.DEBUG`- 用于调试消息（最低严重性）

## 4.1.3 如何记录消息

至于如何使用 scrapy logging 来记录的示例就不再展示了，具体使用方法请看文档：[scrapy logging 使用说明](#)，本库更推荐在调试阶段使用 loguru 来打印日志，会更快捷和明显地查看自己注意的部分。

ayugespidertools 会在 startproject 后默认再 settings 中添加一个日志配置，用于当前项目全局使用，可以在项目的各个目录文件中使用。

以下是如何使用 loguru 的 WARNING 级别记录消息的快速示例：

```
# project_name 为当前所在的 scrapy 项目名称
from <project_name>.settings import logger
logger.warning("This is a warning")
```

因为 Loguru 的理念是：**只有一个 logger**。将日志消息分派给已配置处理程序的对象。Logger 是 loguru 的核心对象，每个日志配置和使用都要通过对其中一个方法的调用。只有一个记录器，因此在使用之前不需要检索一个记录器。具体请查看文档：[loguru 使用说明](#)

所以，你也可以直接使用如下方式，也会在全局中使用同一个 loguru。

```
from loguru import logger

logger.info("this is a info log")
```

## 4.1.4 从 spider 记录

以下是本库 slog 日志在 spider 中的使用示例：

```
import ayugespidertools
```

(continues on next page)

(continued from previous page)

```
class MySpider(ayugespidertools.AyuSpider):

    name = 'myspider'
    start_urls = ['https://scrapy.org']

    def parse(self, response):
        # 此条 (error 级别以下的) 日志默认下只会在控制台输出
        self.slog.info(f"info: Parse function called on {response.url}")
        # 此条日志在默认下会持久化存储至 error.log 中
        self.slog.error(f"error: Parse function called on {response.url}")
```

注：不影响 scrapy 自带的日志记录，可自行选择或同时使用。

### **Logging**

在 ayugespidertools 上学习如何使用日志。



## 扩展 SCRAPY

## 5.1 downloader-middleware

介绍本库中自带的常用 DOWNLOADER\_MIDDLEWARES 中间件。

需要使用本库中的配置，需要在 spider 中修改如下，此为前提：

```
from ayugespidertools.spiders import AyuSpider

# 当前 spider 要继承 AyuSpider
class DemoOneSpider(AyuSpider):
    ...
```

### 5.1.1 1. 随机 UA

使用 fake\_useragent 库中的 ua 信息，在每次发送请求时将随机取 ua 信息，将比较常用的 ua 标识的权重设置高一点，这里是根据 fake\_useragent 库中的打印信息来规划权重的，即类型最多的 ua 其权重也就越高。

#### 1.1. 使用方法

只需激活 DOWNLOADER\_MIDDLEWARES 对应的配置即可。

```
custom_settings = {
    "DOWNLOADER_MIDDLEWARES": {
        # 随机请求头
        "ayugespidertools.middlewares.RandomRequestUaMiddleware": 400,
    },
}
```

若想查看是否正常运行，只需查看其 scrapy 的 debug 日志，或在 spider 中打印 response 信息然后查看其信息即可。

## 5.1.2 2. 代理

### 2.1. 动态隧道代理

本库以快代理为例，其各个代理种类的使用方法大致相同

#### 2.1.1. 使用方法

激活 DOWNLOADER\_MIDDLEWARES 中的动态代理配置。

```
custom_settings = {
    "DOWNLOADER_MIDDLEWARES": {
        # 动态隧道代理激活
        "ayugespidertools.middlewares.DynamicProxyDownloaderMiddleware": 125,
    },
}
```

需要修改此项目中的 VIT 文件夹下的 .conf 对应的配置信息。

```
[KDL_DYNAMIC_PROXY]
PROXY=o668.kdltps.com:15818
USERNAME=***
PASSWORD=***
```

然后即可正常运行。

### 2.2. 独享代理

#### 2.2.1. 使用方法

激活 DOWNLOADER\_MIDDLEWARES 中的独享代理配置。

```
custom_settings = {
    "DOWNLOADER_MIDDLEWARES": {
        # 独享代理激活
        "ayugespidertools.middlewares.ExclusiveProxyDownloaderMiddleware": 125,
    },
}
```

需要修改此项目中的 VIT 文件夹下的 .conf 对应的配置信息。



```
[kdl_exclusive_proxy]
proxy=http://kps.kdlapi.com/api/getkps?orderid=***&num=100&format=json
username=***
password=***
index=1
```

注：index 为在有多个独享代理时取的代理对应的索引值。

### 5.1.3 3. 发送请求库改为 requests

#### 3.1. 使用方法

激活 DOWNLOADER\_MIDDLEWARES 中的对应配置。

```
custom_settings = {
    "DOWNLOADER_MIDDLEWARES": {
        # 替换 scrapy Request 请求为 requests 的中间件
        "ayugespidertools.middlewares.RequestsDownloaderMiddleware": 401,
    },
}
```

然后在 spider 中正常 yield scrapy.Request 即可。

### 5.1.4 4. 发送请求方式改为 aiohttp

#### 4.1. 使用方法

激活 DOWNLOADER\_MIDDLEWARES 中的对应配置。

```
custom_settings = {
    "TWISTED_REACTOR": "twisted.internet.asyncioreactor.AsyncioSelectorReactor",
    "DOWNLOADER_MIDDLEWARES": {
        # 将 scrapy Request 替换为 aiohttp 方式
        "ayugespidertools.middlewares.AiohttpDownloaderMiddleware": 543,
    },
    # scrapy Request 替换为 aiohttp 的配置示例
    "LOCAL_AIOHTTP_CONFIG": {
        "timeout": 2,
        "proxy": "http://127.0.0.1:7890",
        "sleep": 1,
        "retry_times": 3,
    },
}
```

注:

- TWISTED\_REACTOR 的配置在本库的 settings 中就默认打开的, 这里配置是为了演示, 不用再次配置的;
- LOCAL\_AIOHTTP\_CONFIG 为 aiohttp 的全局配置, 一般可用来配置 proxy, timeout, retry\_times, sleep 等全局的参数的;

然后需要构造 AioHttpRequest 或 AioHttpRequestFormRequest 请求对象, 具体的 aiohttp 参数在 args 中配置:

具体使用习惯和方式看个人选择, 如下:

```
# 方式一: 通过 args 参数传值
yield AioHttpRequest(
    url="http://httpbin.org/get?get_args=1",
    callback=self.parse_get_fir,
    meta={
        "meta_data": "这是用来测试 parse_get_fir meta 的功能",
    },
    args=AioHttpRequestArgs(
        method="GET",
        headers={
            "Cookie": "headers_ck_key1=ck; headers_ck_key2=ck",
        },
        cookies={
            "ck_key": "ck",
        },
    ),
    dont_filter=True,
)

# 方式二: 使用 scrapy 传统方式传值
yield AioHttpRequest(
    url="http://httpbin.org/get?get_args=1",
    callback=self.parse_get_fir,
    headers={
        "Cookie": "headers_ck_key1=ck; headers_ck_key2=ck",
    },
    cookies={
        "ck_key": "ck",
    },
    meta={
        "meta_data": "这是用来测试 parse_get_fir meta 的功能",
    },
    dont_filter=True,
```

(continues on next page)

(continued from previous page)

```
)

# 同样, 发送 post 也可以选择两种方式
# 测试 POST 请求示例一 - normal
post_data = {"post_key1": "post_value1", "post_key2": "post_value2"}
yield AioHttpRequest (
    url="http://httpbin.org/post",
    method="POST",
    callback=self.parse_post_fir,
    headers={
        "Cookie": "headers_ck_key=ck; headers_ck_key2=ck",
    },
    body=json.dumps(post_data),
    cookies={
        "ck_key": "ck",
    },
    meta={
        "meta_data": "这是用来测试 parse_post_fir meta 的功能",
    },
    cb_kwargs={
        "request_name": "normal_post1",
    },
    dont_filter=True,
)

# 测试 POST 请求示例一 - aiohttp args
yield AioHttpRequest (
    url="http://httpbin.org/post",
    callback=self.parse_post_fir,
    args=AioHttpRequestArgs (
        method="POST",
        headers={
            "Cookie": "headers_ck_key=ck; headers_ck_key2=ck",
        },
        cookies={
            "ck_key": "ck",
        },
        data=json.dumps(post_data),
    ),
    meta={
        "meta_data": "这是用来测试 parse_post_fir meta 的功能",
    },
    cb_kwargs={
        "request_name": "aiohttp_post1",
```

(continues on next page)

```
    },
    dont_filter=True,
)

# 测试 POST 请求示例二 - normal
yield AiohttpFormRequest (
    url="http://httpbin.org/post",
    headers={
        "Cookie": "headers_ck_key=ck; headers_ck_key2=ck",
    },
    cookies={
        "ck_key": "ck",
    },
    formdata=post_data,
    callback=self.parse_post_sec,
    meta={
        "meta_data": "这是用来测试 parse_post_sec meta 的功能",
    },
    cb_kwargs={
        "request_name": "normal_post2",
    },
    dont_filter=True,
)

# 测试 POST 请求示例二 - aiohttp args
yield AiohttpFormRequest (
    url="http://httpbin.org/post",
    callback=self.parse_post_sec,
    args=AiohttpRequestArgs (
        method="POST",
        headers={
            "Cookie": "headers_ck_key=ck; headers_ck_key2=ck",
        },
        cookies={
            "ck_key": "ck",
        },
        data=post_data,
    ),
    meta={
        "meta_data": "这是用来测试 parse_post_sec meta 的功能",
    },
    cb_kwargs={
        "request_name": "aiohttp_post2",
    },
    ),
```

(continues on next page)

(continued from previous page)

```
dont_filter=True,  
)
```

## 5.2 pipelines

介绍本库中自带的常用 pipelines 管道。

需要使用本库中的配置，需要在 spider 中修改如下，此为前提：

```
from ayugespidertools.spiders import AyuSpider  
  
# 当前 spider 要继承 AyuSpider  
class DemoOneSpider(AyuSpider):  
    ...
```

### 5.2.1 1. Mysql 存储

#### 1.1. 普通存储

AyuFtyMysqlPipeline 为 mysql 存储的普通模式，具有自动创建所需数据库，数据表，自动动态管理 table 字段，表注释，也会自动处理常见（字段编码，Data too long，存储字段不存在等等）的存储问题。

#### 1.1.1 使用方法

只需激活 DOWNLOADER\_MIDDLEWARES 对应的配置即可。

```
custom_settings = {  
    "ITEM_PIPELINES": {  
        # 激活此项则数据会存储至 Mysql  
        "ayugespidertools.pipelines.AyuFtyMysqlPipeline": 300,  
    },  
}
```

然后在 spider 中按照约定的格式进行 yield item 即可，具体请查看 [yield item](#)，然后不用再去管 pipelines 了。

## 1.2. 异步存储

### 1.2.1. twisted 实现

使用 twisted 的 adbapi 实现 Mysql 存储场景下的异步操作

#### 1.2.1.1. 使用方法

同样地，只需激活 DOWNLOADER\_MIDDLEWARES 对应的配置即可。

```
custom_settings = {
    "ITEM_PIPELINES": {
        # 激活此项则数据会存储至 Mysql
        "ayugespidertools.pipelines.AyuTwistedMysqlPipeline": 300,
    },
}
```

## 1.3. 运行日志记录

打开 RECORD\_LOG\_TO\_MYSQL 参数会记录 spider 的运行情况和所依赖的数据库下（带有 crawl\_time 字段的）所有表格的当前采集情况统计。

## 5.2.2 2. MongoDB 存储

这里就直接介绍其依赖方法，因为其它配置与上方 Mysql 场景一模一样

### 2.1. 普通存储

```
# 依赖 AyuFtyMongoPipeline
```

### 2.2. 异步存储

#### 2.2.1. twisted 实现

```
# 依赖 AyuTwistedMongoPipeline
```

### 2.2.2. asyncio motor 实现

```
# 依赖 AsyncMongoPipeline
```

#### *downloader-middleware*

了解本库中的下载中间件及使用方法。

#### *pipelines*

了解本库中的管道及使用方法。





## 构建你的专属库

### 6.1 How-To-Build-Your-Own-Library

本库由 `poetry` 包管理工具构建，任何修改本库后自定义打包等需求请以 `poetry` 官方文档为准。

#### 6.1.1 前言

本库是把 `Scrapy` 的一些常用功能（扩展功能和开发中常用方法）封装成了一个库，方便大家快速使用。

但在使用本库的过程中，你可能会遇到一些问题：

- 比如模板中某些配置不符合你的项目需求；
- 不喜欢项目结构的设计；
- 依赖库的版本不适合你的需求。

像这些可能包含非常个性化的定制，无法适配所有人的喜好，也无法通过 `Pull Requests` 合并来优雅地解决此类问题，这时候你可能会想要修改一些东西，那么你可以参考本文档，来快速构建你自己的专属库。

#### 6.1.2 构建方法

你可以 `clone` 源码后，修改任意方法，修改完成后 `poetry build` 或 `make build` 即可打包并内部使用。

以更新 `pymongo` 版本为例：

- `clone` 项目并准备开发环境

将项目克隆到本地，创建 `python 3.8.1+` 环境并安装 `poetry`，然后在项目根目录下运行 `poetry install` 安装依赖即可。

- 自定义内容

修改你所关注的部分，比如你的项目场景下可能需要其它的日志配置默认值，或添加其它的项目结构模板，更改库名等。

若需要更新本项目的 `pymongo` 依赖库版本为 `x.x.x`, 那只需 `poetry add pymongo==x.x.x` 安装目标版本即可。

- 重打包

修改完毕并测试可用后, 即可通过 `poetry build` 或 `make` 工具的 `make build` 打包即可使用。

### 6.1.3 补充

若你自定义的方法对大多数人都合适的话, 可以尝试将此功能添加到本项目, 但是在此之前请先发相关的 `ISSUES` 确认可行后再开发和提交对应 `PULL REQUESTS`, 以免浪费了你做出的贡献。

#### *How-To-Build-Your-Own-Library*

如何将本库构建成为你的专属库。

**补充说明**

## 7.1 Release notes

### 7.1.1 AyugeSpiderTools 3.0.1 (2023-05-17)

这是一个 major 版本更新，含有 bug 修复、代码优化等。

#### Deprecation removals

- 删除 ayugespidertools 的 cli 名称 -> 改为 ayuge 来管理。

#### Deprecations

- 无。

#### New features

- 修改 item 实现方式，不再通过将字段都存入 alldata 中即可实现动态设置字段的功能，使用更清晰，且能更方便地使用 ItemLoaders 的功能，具体内容及示例请查看 [readthedocs](#) 上对应内容，具体案例请查看 DemoSpider 项目。

#### Bug fixes

- 修复不会创建表注释的问题。

### Code optimizations

- 修改 `dict_keys_to_lower` 和 `dict_keys_to_upper` 的将字典 `key` 转为大写或小写的功能优化为嵌套字典中所有 `key` 都转为大写或小写。
- 将模板中 `settings.py` 中的配置读取放入库中 `update_settings` 实现, 简化 `settings.py` 文件内容。
- 优化 `Makefile` 功能, 简化清理 `__pycache__` 文件夹的功能。
- 修改部分 `typo` 问题。
- 更新 `readthedocs` 内容, 更新测试文件。

### 7.1.2 AyugeSpiderTools 2.1.0 (2023-05-09)

这是一个主要更改了 `scrapy` 依赖库为 2.9.0 版本, 含有 `bug` 修复。

#### Deprecation removals

- `tox` 去除 `windows` 平台的测试场景。

#### Deprecations

- 下一大版本将删除 `ayugespidertools` 的 `cli` 名称 -> 改为 `ayuge` 来管理。

#### New features

- 本库依赖库 `scrapy` 版本升级为 2.9.0。

#### Bug fixes

- 修复使用 `ayuge` 及 `ayuge -h` 命令时, 未显示当前库版本的问题。

#### Code optimizations

- 无。

### 7.1.3 AyugeSpiderTools 2.0.3 (2023-05-06)

此版本为微小变动。

#### Deprecation removals

- 无

#### Deprecations

- 下一大版本将删除 ayugespidertools 的 cli 名称 -> 改为 ayuge 来管理。

#### New features

- 添加 mongodb 的 asyncio 的示例。

#### Bug fixes

- 无

#### Code optimizations

- readthedocs 的 markdown 解析由 recomcommonmark 改为 myst-parser, 以支持更多的 markdown 语法。

### 7.1.4 AyugeSpiderTools 2.0.1 (2023-04-27)

此版本为大版本更新, 修改了项目结构以统一本库及与 scrapy 结合的代码风格, 也有一些功能完善等。最新功能示例请在 [DemoSpider](#) 或 [readthedocs](#) 中查看。

#### Deprecation removals

- 一些 api 变动:
- 一些参数配置变动:

注: 所有配置的 key 都统一改为小写

- 一些使用方法更改:
  - 使用 AioHttpRequest 构造请求时, 由之前的 meta 中的 aiohttp\_args 配置参数, 改为由 args 的新增参数取代, 其参数类型同样为 dict, 也可以为 AioHttpRequestArgs 类型, 更容易输入。

### Deprecations

- 下一大版本将删除 `ayugespidertools` 的 `cli` 名称 -> 改为 `ayuge` 来管理。

### New features

- 丰富 `aihttp` 请求场景，增加重试，代理，`ssl` 等功能。

### Bug fixes

- 无

### Code optimizations

- 更新测试用例。

## 7.1.5 AyugeSpiderTools 1.1.9 (2023-04-20)

这是一个维护版本，具有次要功能、错误修复和清理。

### Deprecation removals

- 无

### Deprecations

- 无

### New features

- 增加 `ayuge startproject` 命令支持 `project_dir` 参数。

```
# 这将在 project_dir 目录下创建一个 Scrapy 项目。如果未指定 project_dir，则 project_dir
→ 将与 myproject 相同。
ayuge startproject myproject [project_dir]
```

## Bug fixes

- 修复模板中 `settings` 的 `CONSUL` 配置信息没有更新为 `v1.1.6` 版本推荐的使用方法的问题。(releases [ayugespidertools-1.1.6](#))
- 修复在 `startproject` 创建项目时生成的 `run.sh` 中的路径信息错误问题。

## Code optimizations

- 添加测试用例。
- 以后的版本发布说明都会在 [ayugespidertools readthedocs](#) 上展示。

## Release notes

查看最近的 AyugeSpiderTools 版本中有哪些变化。