

“网络爬虫”未完待续...

WS00



嵩天

www.python123.org

The Website is the API ...



Requests

自动爬取HTML页面
自动网络请求提交

robots.txt

网络爬虫排除标准



Beautiful Soup

解析HTML页面



Re

正则表达式详解
提取页面关键信息



Scrapy*

*网络爬虫原理介绍
*专业爬虫框架介绍



Projects

实战项目A/B

掌握定向网络数据爬取和网页解析的基本能力

python
弹指之间 · 享受创新

Python网络爬虫与信息提取

04X -Tian

Scrapy爬虫的地位

Python语言最好的爬虫框架

具备企业级专业爬虫的扩展性（7x24高可靠性）

千万级URL爬取管理与部署

简单说：Scrapy足以支撑一般商业服务所需的爬虫能力

Scrapy爬虫的应用展望

持续爬取、商业服务、高可靠性

普通价值：

基于Linux，7x24，稳定爬取输出

商业级部署和应用（scrapyd-*）

千万规模内URL爬取、内容分析和存储

Scrapy爬虫的应用展望

只有Scrapy还不够

高阶价值：

基于docker，虚拟化部署

中间件扩展，增加调度和监控

各种反爬取对抗技术

请关注嵩老师的“MOOC进阶课程”...

面向计算机类专业从业者，讲授高规格产品级实战技术

围绕真实项目为教学设计原型

中国大学MOOC平台 (www.icourses.cn/imooc) 在线开放课程



Python提高课程



Python基础课程

基础课程：语言入门和基本使用



Python网络爬虫与信息提取



Python游戏开发入门



Python数据分析与展示



Python云端系统开发入门



Python机器学习应用



Python科学计算三维可视化

提高课程：各专题领域入门和基本使用