



به نام خدا

آموزش وب کاوی با SCRAPY

مؤلف:

محمد رضا شاقوزی



هرگونه چاپ و تکثیر از محتویات این کتاب بدون اجازه کتبی ناشر ممنوع است. متخلفان به موجب قانون حمایت حقوق مؤلفان، مصنفان و هنرمندان تحت پیگرد قانونی قرار می‌گیرند.

◀ عنوان کتاب: آموزش وب کاوی با SCRAPY

◀ مولف: محمدرضا شاقوزی

◀ ناشر: موسسه فرهنگی هنری دیباگران تهران

◀ ویراستار: مهدیه مخبری

◀ صفحه آرای: نازنین نصیری

◀ طراح جلد: داریوش فرسایی

◀ نوبت چاپ: اول

◀ تاریخ نشر: ۱۴۰۱

◀ چاپ و صحافی: درج عقیق

◀ تیراژ: ۱۰۰۰ جلد

◀ قیمت: ۱۰۲۰۰۰۰ ریال

◀ شابک: ۹۷۸-۶۲۲-۲۱۸-۶۴۸-۷

نشانی واحد فروش: تهران، خیابان انقلاب، خیابان دانشگاه

-تقاطع شهدای ژاندارمری-پلاک ۱۵۸ ساختمان دانشگاه-

طبقه دوم-واحد ۴ تلفن ها: ۶۶۹۶۵۷۴۹-۲۲۰۸۵۱۱۱

فروشگاههای اینترنتی دیباگران تهران :

WWW.MFTBOOK.IR

www.dibagaran-tehran.com

سرشناسه: شاقوزی، محمدرضا، ۱۳۷۳-
عنوان و نام پدیدآور: آموزش وب کاوی با SCRAPY
/مولف: محمدرضا شاقوزی؛
ویراستار: مهدیه مخبری.
مشخصات نشر: تهران: دیباگران تهران: ۱۴۰۱
مشخصات ظاهری: ۱۲۲ ص: جدول
شابک: ۹۷۸-۶۲۲-۲۱۸-۶۴۸-۷
وضعیت فهرست نویسی: فیپا
موضوع: داده کاوی Data Mining
موضوع: اسکریپت (software framework) scrapy
موضوع: داده های کلان big data
رده بندی کنگره: ۷۶/۹ QA
رده بندی دیویی: ۰۰۶/۳۱۲
شماره کتابشناسی ملی: ۹۱۲۴۱۱۴

نشانی اینستاگرام دیبا dibagaran_publishing نشانی تلگرام: @mftbook

هر کتاب دیباگران، یک فرصت جدید علمی و شغلی.

هر گوشی همراه، یک فروشگاه کتاب دیباگران تهران.

از طریق سایتهای دیباگران، در هر جای ایران به کتابهای ما دسترسی دارید.

فهرست مطالب

مقدمه ناشر	۵
مقدمه مؤلف	۶
۱ وب اسکریپینگ	۷
۱-۱- چرا وب اسکریپینگ؟	۷
۲ آموزش Beautifulsoup	۱۰
۱-۲- زبان نشانه‌گذاری HTML	۱۰
۲-۲- استفاده از مرورگر به‌عنوان ابزار توسعه	۱۲
۳-۲- Cascading Style Sheets: CSS	۱۷
۴-۲- کتابخانه Beautiful Soup	۲۲
۵-۲- مطالب بیشتر در مورد Beautiful Soup	۳۵
۳ بررسی عمیق‌تر HTTP	۳۹
۱-۳- کار با Form و درخواست POST	۳۹
۲-۳- سایر روش‌های درخواست HTTP	۵۲
۳-۳- اطلاعات بیشتر دربارهٔ هدرها (Headers)	۵۵
۴-۳- قالب‌های محتوا (Binary, Json,...)	۵۸
۴ از خراش وب تا خزش وب	۶۴
۱-۴- خزش وب (web crawling) چیست؟	۶۴
۲-۴- خزش وب در پایتون	۶۷
۳-۴- ذخیره نتایج در یک پایگاه داده	۷۰
۵ علم داده	۸۴
۱-۵- فرایند علم داده	۸۴
۲-۵- جایگاه وب اسکریپینگ در فرایند علم داده	۸۸
۶ مثال‌های کاربردی از اسکریپ بدون استفاده از فریمورک	۹۱
۱-۶- اسکریپ Hacker News	۹۱

۹۳.....۲-۶- اسکریپ مخزن گوگل در گیت‌هاب

۹۵ فریمورک Scrapy

۹۶.....۱-۷- نصب Scrapy

۹۶.....۲-۷- راه‌اندازی یک spider جدید

۹۷.....۳-۷- نوشتن یک اسکریپر ساده

۹۹.....۴-۷- قوانین Spider

۱۰۴.....۵-۷- ساختن Items

۱۰۷.....۶-۷- خروجی دادن Items

۱۰۸.....۷-۷- خط لوله Items

۱۱۳.....۸-۷- لاگ گرفتن با Scrapy

۱۱۴.....۹-۷- تعریف Encoding برای متون فارسی

۱۱۵ نمونه پروژه‌های Scrapy

۱۱۵.....۱-۸- کرول سایت نوبت دات آی آر؛ سایت نوبت‌دهی پزشکان

۱۱۷.....۲-۸- کرول سایت دیوار؛ آگهی‌های خودرو

۱۲۰.....۳-۸- کرول اخبار سایت ایسنا

خط‌مشی انتشارات مؤسسه فرهنگی هنری دیباگران تهران در عرصه کتاب‌هایی با کیفیت عالی است که بتواند
خواسته‌های به روز جامعه فرهنگی و علمی کشور را تا حد امکان پوشش دهد.
هر کتاب دیباگران تهران، یک فرصت جدید شغلی و علمی

حمد و سپاس ایزد منان را که با الطاف بی‌کران خود این توفیق را به ما ارزانی داشت تا بتوانیم در راه ارتقای دانش عمومی و فرهنگی این مرز و بوم در زمینه چاپ و نشر کتب علمی و آموزشی گام‌هایی هرچند کوچک برداشته و در انجام رسالتی که بر عهده داریم، مؤثر واقع شویم.

گسترده‌گی علوم و سرعت توسعه روزافزون آن، شرایطی را به وجود آورده که هر روز شاهد تحولات اساسی چشمگیری در سطح جهان هستیم. این گسترش و توسعه، نیاز به منابع مختلف از جمله کتاب را به عنوان قدیمی‌ترین و راحت‌ترین راه دستیابی به اطلاعات و اطلاع‌رسانی، بیش از پیش برجسته نموده است.

در این راستا، واحد انتشارات مؤسسه فرهنگی هنری دیباگران تهران با همکاری اساتید، مؤلفان، مترجمان، متخصصان، پژوهشگران و محققان در زمینه‌های گوناگون و مورد نیاز جامعه تلاش نموده برای رفع کمبودها و نیازهای موجود، منابعی پُر بار، معتبر و با کیفیت مناسب در اختیار علاقمندان قرار دهد.

کتابی که در دست‌دارید تألیف "جناب آقای محمدرضا شاقوزی" است که با تلاش همکاران ما در نشر دیباگران تهران منتشر گشته و شایسته است از یکایک این گرامیان تشکر و قدردانی کنیم.

با نظرات خود مشوق و راهنمای ما باشید

با ارائه نظرات و پیشنهادات و خواسته‌های خود، به ما کمک کنید تا بهتر و دقیق‌تر در جهت رفع نیازهای علمی و آموزشی کشورمان قدم برداریم. برای رساندن پیام‌هایتان به ما از رسانه‌های دیباگران تهران شامل سایتهای فروشگاهی و صفحه اینستاگرام و شماره‌های تماس که در صفحه شناسنامه کتاب آمده استفاده نمایید.

مدیر انتشارات

مؤسسه فرهنگی هنری دیباگران تهران
dibagaran@mftplus.com

مقدمه

در قرن حاضر با ظهور تکنولوژی‌های ارتباطات، شاهد انفجار اطلاعات هستیم. کمتر موضوعی پیدا می‌شود که برای آن محتواهای بسیاری بر روی اینترنت نتوان پیدا کرد. از طرفی می‌دانیم هرچه اطلاعات بیشتری داشته باشیم در تصمیم‌گیری‌های روزمره سود بیشتری نصیبمان می‌شود. فرض کنید می‌خواهید یک ماشین کارکرده بخرید. در حالت معمولی، شما از طریق دوستان و آشنایانتان از قیمت خودرو و خریدار مطلع می‌شوید. اما اگر بتوانید به صورت هدفمند از یک سایت خرید و فروش دست دوم، قیمت خودروی مورد نظر را در طول روزهای مختلف جمع‌آوری کنید، چه مقدار سود خواهید کرد؟ حتی شاید بتوانید با یک دقت مناسبی، قیمت فردای آن خودرو را نیز پیش‌بینی کنید. یا شاید به عنوان یک مشاور معتمد برای دوستان و آشنایانتان عمل کنید و از ضرر احتمالی آن‌ها جلوگیری کنید. **وب کاوی با اسکریپی** شما را برای ساخت چنین ابزارهایی راهنمایی می‌کند.

این کتاب از مباحث پایه شروع می‌کند و به پروژه‌های تقریباً کاربردی ختم می‌شود. برای راحتی ارتباط خواننده، بین فصل‌هایی که به مفاهیم نظری می‌پردازند، یک فصل کاربردی قرار داده شده است تا روند مطالعه کتاب فرسایشی نشود.

لازم به ذکر است که کدهای پروژه و مراحل انجام آن در مخزن گیت‌هاب مربوط به این کتاب در لینک زیر قابل دسترس است. پس اگر کیفیت چاپ کدها یا لینک‌هایی که به صورت `precent encoding` در متن نمایش داده شده‌اند (af/۸c%da/)، حس خوبی به شما نداد، چیزی را از دست نداده‌اید. زیرا می‌توانید تمام کدها را از طریق لینک زیر به تفکیک فصل مشاهده کنید:

https://github.com/softrebel/web_mining_with_scrapy

البته برای حداکثر یادگیری پیشنهاد می‌شود که به هیچ عنوان کدها را کپی نکنید و خودتان کدنویسی کنید. همچنین به علت کمبود فضا، بعضی از خطوط کد پایتون در کتاب با کاراکتر "\" به دو خط تقسیم شده و همچنان قابل اجرا است. این در حالی است که در مخزن سعی شده است کدها به کاربرد نزدیک‌تر باشد. لذا این تفاوت را به منزله اشتباه بودن کدهای کتاب در نظر نگیرید.

به دلیل اینکه وب‌سایت‌های فعلی ایستا نیستند، لازم شد بعضی از تمرین‌ها و مثال‌ها را بر روی یک سایت مختص این کتاب پیاده‌سازی کنم که در طول زمان تغییری ایجاد نشود. دامنه زیر برای این کار تهیه شده است:

<http://www.pycrawling.ir/>

البته این ایراد برای پروژه‌ها نیز وارد است و ممکن است در سال آینده ساختار آن‌ها تغییر کند. در صورت تغییرات، در ویرایش‌های بعدی و در مخزن گیت، آن‌ها را لحاظ می‌کنیم. همچنین شما می‌توانید انتقادات و پیشنهادات خود را با رایانامه sh.mohammad66@gmail.com با بنده در میان بگذارید تا در ویرایش‌های بعدی در نظر گرفته شود.

در انتها لازم است از صبر و زحمات تمام دست‌اندرکاران مجموعه دیباگران تشکر نمایم.

ارادتمند

محمد رضا شاقوزی

