# باسمه تعالی درس: بازیابی اطلاعات



نام و نام خانوادگی: تینا توکلی و هادی امینی

شماره دانشجویی: 9912762270 و 9912762270

شماره تمرین: 03

# توضيح كد:

از كتابخانه selenium براى استخراج داده از وب سايت IEEE Xplore استفاده مي كنيم.

# متغير سراسري:

driver: یک متغیر سراسری برای نگهداری از نمونه WebDriver. در ابتدا به None تنظیم می شود.

# توابع:

load\_homepage:()

یک نمونه جدید Chrome WebDriver ایجاد می کند.

به URL صفحه اصلي IEEE Xplore هدايت مي شود.

نمونه WebDriver را برمی گرداند.

search(query):

منتظر می ماند تا فیلد و رودی جستجو قابل کلیک باشد.

کادر جستجو را پیدا می کند و query ارائه شده را وارد می کند.

دكمه جستجو را بيدا مي كند و روي آن كليك مي كند.

get\_papers:()

یک لیست خالی paper\_list برای ذخیره لینک های مقاله ایجاد می کند.

5 بار حلقه مي زند (قابل تنظيم):

منتظر حضور نتایج جستجو می ماند.

تمام عناصر با کلاس result-item-align را پیدا می کند.

هر مورد نتیجه را پیمایش می کند:

publisher-info-container را پیدا می کند.

نوع مقاله را با استفاده از نام های تگ و متن از دومین عنصر span استخراج می کند.

اگر نوع "مقاله كنفرانس" باشد:

لینک مقاله را با استفاده از ویژگی href آن استخراج می کند.

لینک را به paper list اضافه می کند.

next\_page) را برای کلیک بر روی دکمه "صفحه بعد" (در صورت وجود) فراخوانی می کند.

paper\_list را برمی گرداند.

next\_page:()

ظرف دكمه "صفحه بعد" (next-page-set) را بيدا مي كند.

منتظر می ماند تا دکمه قابل کلیک باشد.

از اجرای جاوا اسکریپت (execute\_script) برای کلیک بر روی دکمه استفاده می کند (مشکلات احتمالی تعامل را مدیریت می کند).

### select\_value(string):

رشته ورودی را با استفاده از جداکننده ": " تقسیم می کند و مقدار بعد از دو نقطه را برمی گرداند (فرمت خاصی برای استخراج داده فرض می کند).

extract\_data(paper\_link):

به paper\_link ارائه شده هدایت می شود.

منتظر مى ماند تا عنصر عنوان مقاله موجود باشد.

یک دیکشنری خالی info برای ذخیره داده های استخراج شده برای این مقاله ایجاد می کند.

اطلاعات مختلف را با استفاده از عبارات XPath، نام کلاس ها و انتخاب کننده های CSS خاص استخراج می کند (برخی در بلوک های try...except برای مدیریت خطاهای احتمالی پیچیده شده اند):

# سلكتورها:

### :title .1

- انتخاب کننده: By.CLASS\_NAME, 'document-title'
- توضیح: این سلکتور به دنبال عنصری با نام کلاس document-title میگردد. این نام کلاس احتمالاً به عنصری اشاره میکند که عنوان سند را در صفحه وب نشان میدهد.

#### :pages .2

- انتخاب كننده: [[(', '/div[@class='u-pb-1'][strong[contains(text(), 'Page(s):')]]
  - توضیح: این عبارت XPath یک ساختار عنصری خاص را هدف قرار میدهد:
- //(div @class='u-pb-1: عناصر div و ا با نام کلاس u-pb-1 پیدا میکند. این کلاس ممکن است به بخشی حاوی اطلاعات صفحه مربوط باشد.

• [[(':strong[contains(text(), 'Page(s)]: عناصر div را فقط به مواردی که شامل یک عنصر strong با متنی حاوی "(Page(s):" هستند، فیلتر میکند. این نشان میدهد که عنصر strong برچسب شماره صفحه را نگه میدارد.

#### : Cites in Papers.3

- انتخاب کننده: -By.XPATH, "//button[div[text()='Papers']]//div[@class='document-banner-metric انتخاب کننده: -count']
  - توضیح: این عبارت XPath یک عنصر دکمه را هدف قرار میدهد که به Cites in Papers مربوط می شود:
- //[button[div[text()='Papers"]: دکمه هایی را پیدا میکند که شامل یک عنصر div با متن "Papers" هستند. این نشان می دهد که دکمه به Cites in Papers مربوط می شود.
- //div[@class='document-banner-metric-count]: در داخل آن دکمهها، به دنبال یک عنصر div فنصر div استنادات را نگه کلاس احتمالاً تعداد واقعی استنادات را نگه می درد.

#### : Cites in Patent.4

- انتخاب کننده: By.XPATH, "//div[contains(@class, 'document-banner-metriccontainer')]//button[div[text()='Cites in'][following-"sibling::div[text()='Patent']]]/div[@class='document-banner-metric-count']
- توضیح: این عبارت XPath پیچیده یک عنصر دکمه را هدف قرار میدهد که به Cites in Patent مربوط می شود:
- //[contains(@class, 'document-banner-metric-container']: یک عنصر div را پیدا میکند که نام div است. این کلاس آن شامل "document-banner-metric-container" است. این کلاس ممکن است به یک ظرف بر ای معیار های سند مر بوط باشد.
  - //[[[[following-sibling::div[text()='Patent]: در داخل آن ظرف، به دنبال دکمههایی با یک div حاوی "Cites in" میگردد که به دنبال آن یک div همتراز با "Patent" وجود دنبال دکمههایی با یک Cites in Patent" وجود دارد. این ساختار نشان میدهد که دکمه به طور خاص برای Cites in Patent است.
  - /div[@class='document-banner-metric-count']: در نهایت، عنصر div را در داخل دکمه انتخاب میکند که احتمالاً تعداد Cites را نگه میدارد.

# : Full Text Views .5

- By.XPATH, "//button[div[div[text()='Full']]]/div[@class='document-banner-metric- انتخاب کننده: -count']
  - توضیح: این عبارت XPath یک عنصر دکمه را پیدا میکند که به باز دیدهای متن کامل مربوط می شود:
- /div @class='document-banner-metric-count': عنصر div و در داخل دکمه انتخاب میکند که احتمالاً تعداد باز دیدهای متن کامل را نگه میدار د.

#### : Publisher .6

• انتخاب كننده : ("By.XPATH, '//span[contains(By.CLASS NAME,"publisher-info-container")

• توضیح (تلاش 2): این عبارت XPath جایگزین دیگری را ارائه میدهد که یک ساختار span تو در تو را در داخل عنصری با کلاس حاوی "publisher-info-container" هدف قرار میدهد.

#### : DOI .7

- انتخاب کننده: By.CLASS\_NAME, 'stats-document-abstract-doi'
- توضیح: این سلکتور به دنبال عنصری با نام کلاس stats-document-abstract-doi میگردد. این نام کلاس احتمالاً به عنصری اشاره میکند که شناسه شیء دیجیتال (DOI) را در خود جای داده است.

#### : Date of Publication .8

- انتخاب کننده: By.CLASS NAME, 'doc-abstract-confdate'
- توضیح: این سلکتور به دنبال عنصری با نام کلاس doc-abstract-confdate میگردد. این نام کلاس احتمالاً به عنصری اشاره میکند که Date of Publication را در خود جای داده است.

#### : Abstract .9

- "[g-0"]//div[@xplmathjax row By.XPATH, '//div[@class="abstract-text انتخاب كننده:
- توضیح: این عبارت XPath یک عنصر div را با کلاس g-0 row abstract-text و حاوی یک عنصر div با ویژگی xplmathjax هدف قرار میدهد. این ساختار احتمالاً به عنصری اشاره میکند که Abstract را در خود جای داده است.

#### : Published in .10

- انتخاب کننده: a By.CSS\_SELECTOR, '.stats-document-abstract-publishedIn'
- توضیح: این سلکتور CSS به دنبال عنصر a با کلاس stats-document-abstract-publishedIn میگردد. این عنصر احتمالاً پیوند به نشریه ای را که سند در آن منتشر شده است، در خود جای داده است.

#### : Authors .11

# • توضيح:

- این کد ابتدا عنصر دکمه با شناسه authors را پیدا میکند و بر روی آن کلیک مه،کند.
  - سیس تمام عناصر div با کلاس author-card را پیدا میکند.
- برای هر عنصر author-card، متن آن را به یک دیکشنری با کلیدهای name و from (در صورت وجود)
  تبدیل میکند.
  - در نهایت، لیست دیکشنریهای نویسنده را در info ذخیره میکند.

#### :IEEE Keywords .12

#### • توضيح:

- این کد ابتدا عنصر دکمه با شناسه keywords را پیدا میکند و بر روی آن کلیک میکند.
- سپس تمام عناصر a را با داده tealium\_data حاوی "IEEE Keywords" پیدا میکند.
- در نهایت، لیست متن عناصر a را در info با عنوان Keywords IEEE ذخیره میکند.

#### :Author Keywords .13

### و توضيح:

- این کد ابتدا عنصر دکمه با شناسه keywords را پیدا میکند و بر روی آن کلیک میکند.
- سپس تمام عناصر a را با داده tealium\_data حاوی "Author Keywords" پیدا میکند.
- در نهایت، لیست متن عناصر a را در info با عنوان Keywords Author ذخیره میکند.

و در انتها دیکشنری info حاوی داده های استخراج شده را برمی گرداند و آنها را با ()save\_json ذخیره میکند .

### چالش ها:

# 1. در بخش Author Keywords : لیست خالی Author Keywords

با استفاده از XPath به دنبال عناصری با data-tealium\_data حاوی "Author Keywords" هستیم و سپس متن آنها را در لیست Author Keywords ذخیره میکنیم. ولی در برخی موارد ممکن است Bala اضافی را انتخاب میکرد که دارای هیچ متنی نیستند .

# راه حل:

با استفاده از شرط (len (e.text) 0 > این موضوع را هندل کردیم که به معنی فقط در صورتی که عنصر و دارای متن باشد، آن را به لیست اضافه میکنیم.

### 2. کلیک Newest:

تأخیر قابل توجهی در کلیک روی دکمه "Newest" برای مرتب سازی نتایج وجود دارد.

#### علت:

دلایل احتمالی متعددی برای این تأخیر وجود دارد:

- 1. بارگذاری ناهمزمان: ممکن است عناصر صفحه وب به طور ناهمزمان بارگذاری شوند، به این معنی که دکمه "Newest" ممکن است بلافاصله پس از بارگذاری صفحه قابل کلیک نباشد.
  - 2. پیچیدگی صفحه: اگر صفحه وب از نظر ساختاری پیچیده باشد، ممکن است یافتن دکمه "Newest" برای Selenium و مان بیشتری طول بکشد.
  - 3. تأخیر شبکه: اگر اتصال شبکه ضعیف باشد، ممکن است ارسال در خواست کلیک به دکمه "Newest" به سرور زمان بیشتری طول بکشد.

# راه حل:

1. استفاده از WebDriverWait: از کتابخانه WebDriverWait: از کتابخانه WebDriverWait: ارای صبر کردن تا دکمه "Newest" قابل کلیک شود قبل از کلیک بر روی آن استفاده شده است. این کار با استفاده از روش until انجام می شود که

- یک شرط را به عنوان ورودی می پذیرد. در این مورد، شرط این است که دکمه با استفاده از EC.element to be clickable
- 2. استفاده از try...except: در نهایت، از یک بلوک try...except برای تلاش برای کلیک بر روی دکمه "Newest" با استفاده از روش click استفاده شده است. اگر این کار به دلیل ElementNotInteractableException شکست بخورد، از driver.execute\_script برای کلیک بر روی دکمه با استفاده از JavaScript استفاده می شود. این روش می تواند در مواردی مفید باشد که Selenium قادر به کلیک بر روی دکمه با استفاده از روش های معمول نیست.

# 3. وجود نداشتن برخی موارد در برخی مقالات:

برخی از صفحات وب ممکن است برخی عناصر را نداشته باشند. این امر می تواند منجر به خطا در هنگام تلاش برای یافتن این عنصر شود.

### راه حل:

برای حل این مشکل، از یک بلوک try...except استفاده شده است. در داخل بلوک try، کد تلاش می کند عناصر را با استفاده از selector ارائه شده پیدا کند و مقدار آنرا تنظیم کند. اگر این کار به دلیل خطا (مانند (ElementNotFoundException ننظیم می شود.

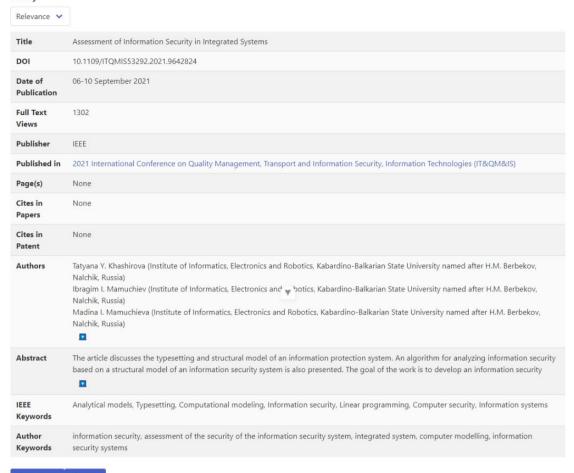
### خروجي:

شامل دو فایل با پسوند json است که یکی برای relevance و newest است که داخل هر کدام شامل تعدادی مقاله کنفرانسی هست.

# بخش اضافي:

یک رابط کاربری تحت وب برای نمایش محتویات فایلهای خروجی ایجاد کردیم.

#### Sort By





Paper 4 of 110

#### مشاركت:

10	امینی
10	توكلى