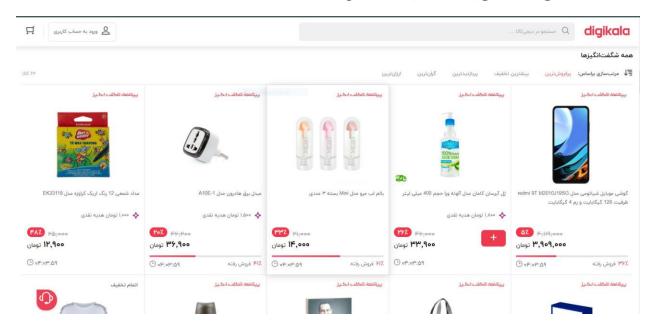
#### شرح پروژه

هدف از پروژه web scraping صفحه فروش شگفت انگیز دیجیکالا، سرعت بخشیدن به روند خرید و فروش و کسب آسان تر اطلاعات هست. در این پروژه از زبان برنامه نویسی پایتون و پکیج های pandas, beautifulSoup4, urllib بهره گرفته شده است. ابتدا به بررسی صفحه اصلی پرداخته و سپس کد توضیح داده میشود.



### https://www.digikala.com/incredible-offers/

در این صفحه قسمت های متعددی وجود دارند که بخش بالا، قسمت مد نظر ما هست. که همه کالا ها با تخفیف مذکور لیست شده اند. در سورس صفحه در این قسمت اطلاعات زیر موجود هستند:



که اگر مرتب شوند دیده میشود که داده ها در تگ های متناسب قسمت بندی شده اند.

نام ها در تگ div class={c-product-box\_title js-ab-not-app-incredible-product}

قیمت ها در تگ {c-price\_\_value-wrapper js-product-card-price} قیمت ها در تگ

تخفیف ها در تگ {c-price\_\_discount-oval}

با استخراج داده های این تگ ها داده مورد نیاز خود را بدست میآوریم.

#### در توضیح کد

ابتدا کتابخانه های مورد نیاز را وارد میکنیم

```
from bs4 import BeautifulSoup
from urllib.request import Request, urlopen
import pandas as pd
```

سپس آدرس سایت و برای داده های خود 3 لیست در نظر میگیریم. همچنین برای تعداد صفحات مورد بررسی و تعداد کالا ها متغیر هایی در نظر میگیریم

```
# https://www.digikala.com/incredible-offers/
originalUrl = "https://www.digikala.com"
urlLink = originalUrl + "/incredible-offers/"

names = []
offs = []
prices = []
noPages = 3
itemCount = 0
```

وارد حلقه برای تعداد صفحات میشویم، برای هر صفحه ابتدا لینک کامل را ساخته، سپس به آن سایت و سرور درخواست فرستاده و جواب را میخوانیم و به کمک کتابخانه bs4 اطلاعات صفحه را پردازش میکنیم

```
for i in range(1, noPages):
    url = urlLink + '?pageno=' + str(i)

    r = Request(url, headers={'User-Agent':'Mozilla/5.0'})
    webpage = urlopen(r).read()
    soup = BeautifulSoup(webpage, 'html.parser')
```

بعد از دریافت page source، مطابق تگ های توضیح داده شده در قسمت قبل، بخش کالا های خود را جدا کرده و در آن بخش ها داده های مورد نیاز خود را جدا میکنیم، همچنین به ازای هر کالا مقدار متغیر شمارش کالا یکی اضافه میشود

```
for row in soup.findAll('div', attrs={'class':'c-product-list_item js-product-list-content'}):
    itemCount += 1
    rName = row.find('div', attrs={'class':'c-product-box_title js-ab-not-app-incredible-product'})
    rPrice = row.find('div', attrs={'class':'c-price_value-wrapper js-product-card-price'})
    rOff = row.find('div', attrs={'class':'c-price_discount-oval'})
```

داده اسم کالا در سورس صفحه شامل فاصله های خالی و اضافی هست. این مشکل در تابع replaceMultiple رفع شده است، بدین صورت که یک dict برای تعویض فاصله های اضافی با فضای خالی در نظر گرفته شده، که در این تابع به ترتیب این تعویض ها انجام میشوند

```
def replaceMultiple(txt, dic):
    for i, j in dic.items():
        txt = txt.replace(i, j)
    return txt
```

همچنین در صورت اتمام تخفیف، داده قیمت مقدار None به خود میگیرد که برای رفع آن مشکل از تابع None داده شوند استفاده شده. در این تابع شرطی بررسی میشود که اگر کالا موجود بود، قیمت آن و اگر نبود خروجی None داده شوند

```
def checkActive(ro):
    if ro != None:
        return ro.text
    else:
        return None
```

برای راحتی نوشتار، این دو عملیات در تابع result باهم ادغام شده اند

```
def result(inp):
    global replaceDict
    out1 = checkActive(inp)
    if out1:
        return replaceMultiple(out1, replaceDict)
    else:
        return 'not available'
```

داده تخفیف نیز در div مخصوص خود، در قسمت span قرار دارد. با توجه به مشکلات استخراج داده span، برای بدست آوردن مقدار تخفیف، از عملیات روی متن استفاده شده است، بدین صورت که متن بین تگ های span و span/ خارج میشود. این عملیات در تابع extractText ییاده سازی شده است

```
def extractText(txt, sth):
    tagStart = '<'+sth+'>'
    tagEnd = '</'+sth+'>'
    res1 = txt.find(tagStart)
    res2 = txt.find(tagEnd)
    if (res1 != -1):
        theIt = txt[res1+len(tagStart):res2]
        return theIt
    else:
        return None
```

پس از آماده سازی اولیه داده های نام و قیمت و تخفیف، آنها را در لیست های مربوطه اضافه میکنیم

```
name = result(rName)
price = result(rPrice)
off = extractText(str(rOff), 'span')
names.append(name)
prices.append(price)
offs.append(off)
```

پس از بررسی همه سایت ها(که در این مورد صفحات 1 و 2 موجود بودند) داده ها را در یک dict ذخیره میکنیم و سپس آن را به فرمت DataFrame از پکیج pandas در میآوریم. هدف از این کار راحتی و امکانات کار با داده در pandas هست

```
dic = {'Name':names, 'Price':prices, 'Off':offs}
df = pd.DataFrame(dic)
```

سپس میتوانیم داده ها را به فرمت CSV تبدیل کنیم تا در سایر برنامه ها نیز قابلیت خواندن داشته باشد، لازم به ذکر است چون داده ها به زبان فارسی هستند، از utf-8 encoding استفاده میکنیم

```
df.to_csv('file1.csv', index=False, encoding='utf-8')
```

- داده هایی که اتمام موجودی باشند، قیمت و تخفیف نخواهند داشت
- داده هایی که تخفیف آنها به پایان رسیده باشد، در ستون تخفیف مقداری ندارند.

	Name	Price	Off
0	ظر redmi 9T M2010J19SG گوشی موبایِل شیائومی مدل	ئومان,۳٫۹۰۹	۵٪
1	ڑل آبرسان کامان مدل آلوئه ورا حجم 400 میلی لیئر	ئومان ۹۰۰ ۳۳٫	19%
2	بسته ۳ عددی Mini بالم لب مرو مدل	ئومان.٠٠٠	۳۳٪
3	A10E-1 مبدل برق هادرون مدل	ئومان.٩٠٠ ٣٦	۲۰٪
4	EK33118 مداد شمعي 12 رنگ اريک کراوزه مدل	ئومان ۱۲٫۹۰۰	44/
61	Ma-1 کیف اداری چرم ما مدل	ئومان,۶۳۰	۳۶/
62	BRIX Blocks ساختنی مدل 1-2007	ئومان.٠٠٠٧	۴./
63	كيِف رودوشي گوگانا كد 300301	not available	None
64	SRM90جانماز ترمه طرح شاه عباسی کد	not available	None
65	C101كفش روزمره زنانه شيفر مدل 5330	not available	None

66 rows × 3 columns

## همچنین برای راحتی بیشتر، در قسمت اضافی پروژه، داده های برست آمده در یک صفحه وب نیز ذخیره میشوند

# df.to\_html('file2.html')

	N	ame	Price	Off
0	۱۹۰۰ ظرفیت 128 گیگابایت و رم 4 گیگابایت redmi 9T M2010J19SG گوشی موبایل شیائومی مدل		تومان ۳,۹۰۹,۰۰۰	
$\vdash$	ر المستقامة على المستقامة المستقام المستقام المستقامة المستقامة المستقام المستقام المستقام المست	_		۲۶٪
$\vdash$	بسته ۳ عددی Mini بالم لب مرو مدل			٣٣٪
3	A10E-1 مبدل برق هادرون مدل		تومان ۳۶٫۹۰۰	۲۰٪
4	EK33118 مداد شمعی 12 رنگ اریک کر اوز ه مدل		تومان،۹۰۰	۴۸٪
5	بسته 500 عددی A+ فولیو پر ایم مدل A4 کاغذ		تومان ۲۶٫۵۰۰	19%
6	FCLT3337 ساک لوازم کودک و نوزاد فوروارد کد		تومان۱۶۹٫۰۰۰	٧۴٪
7	کتاب اثر مرکب اثر دارن هار دی انتشار ات بار ان خرد		تومان۱۹٫۶۰۰	۸٠%
8	حجم 150 میلی لیتر Bvlgari اسپری ضد تعریق مردانه پر ستیژ مدل		تومان،۹۰۰	19%
9	تى شرت زنانه كد 120		تومان.٠٠٠	None
10	حجم 85 میلی لینر Warrior ادو پرفیوم مردانه اسکلاره مدل		تومان۱۵۵٫۹۵۰	None
11	کتاب پنج قدم فاصله اثر ریچل لیپینکات نشر میلکان		تومان۳۶٫۸۰۰	۳۲٪
12	مجموعه 6 عددی A211 جور اب مر دانه مدل نویدکد		, , ,	۵٧٪
13	FCLT3060 کیف ابزار فوروارد مدل		,	۵٩٪
	ظرفیت 10000 میلی آمپر ساعت PB01 شار ژر همراه رانیک پرایم مدل		,	۴۵٪
15	2-FCLT402 کوله پشتی 65 لینزی فوروارد مدل	_	, ,	۴١٪
	روان نویس طرح یونیکورن کد 9911		, - ,	۲۸٪
17	سرویس خواب ژینورا مدل هانا یکنفره 2 تکه		تومان۴۰۰,۰۰۰	None
امداا	Ite may be a see		. 1. A 24	wu*/