

Scrapy CrawlSpider Notlar

Terminal komutları ve anlamları:

#genspider içerisindeki şablonları görmek için terminal komutu.

scrapy genspider -l

Available templates:

basic

crawl

csvfeed

xmlfeed

#şablon oluşturmak:

scrapy genspider deneme x

```
import scrapy
```

```
class DenemeSpider(scrapy.Spider):  
    name = 'deneme'  
    allowed_domains = ['x']  
    start_urls = ['http://x/']
```

```
def parse(self, response):  
    pass
```

#Bu ve aşağıdaki kodun base class farklarına dikkat edin lütfen. Yazdığınız kodu çalıştırmak için iki ayrı komut kullanacaksınız.

#buradaki “deneme” spider adını, x ise başlangıç url’ini temsil eder.

scrapy genspider -t crawl trendy http://www.trendyol.com/

#-t template(şablon) temsil eder, crawl seçtiğimiz şablonu, trendy ise spider adını.

En son ise başlangıç Url’ini yazıyoruz.

```

import scrapy
from scrapy.linkextractors import LinkExtractor
from scrapy.spiders import CrawlSpider, Rule

class TrendySpider(CrawlSpider):
    name = 'trendy'
    allowed_domains = ['www.trendyol.com']
    start_urls = ['http://www.trendyol.com/']

    rules = (
        Rule(LinkExtractor(allow=r'Items/'), callback='parse_item', follow=True),
    )

    def parse_item(self, response):
        item = {}
        #item['domain_id'] = response.xpath('//input[@id="sid"]/@value').get()
        #item['name'] = response.xpath('//div[@id="name"]').get()
        #item['description'] = response.xpath('//div[@id="description"]').get()
        return item

```

Allowed domain kısmını düzeltmemiz gerektiğinin farkına varmamız gerekiyor.

`allowed_domains = ['www.trendyol.com']` \Rightarrow yanlış kullanım.

`allowed_domains = ['trendyol.com']` \Rightarrow doğru kullanım.

Allowed domain kısmı bizim için çok önemli çünkü spider'ımızın domainin dışına çıkmasını istemiyoruz.

Rules

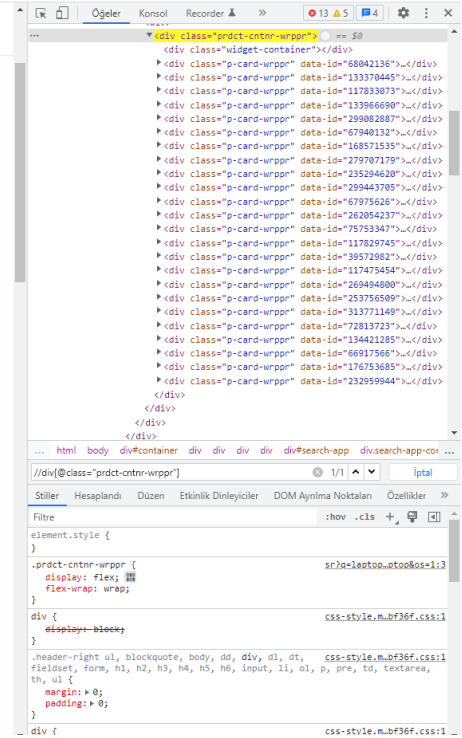
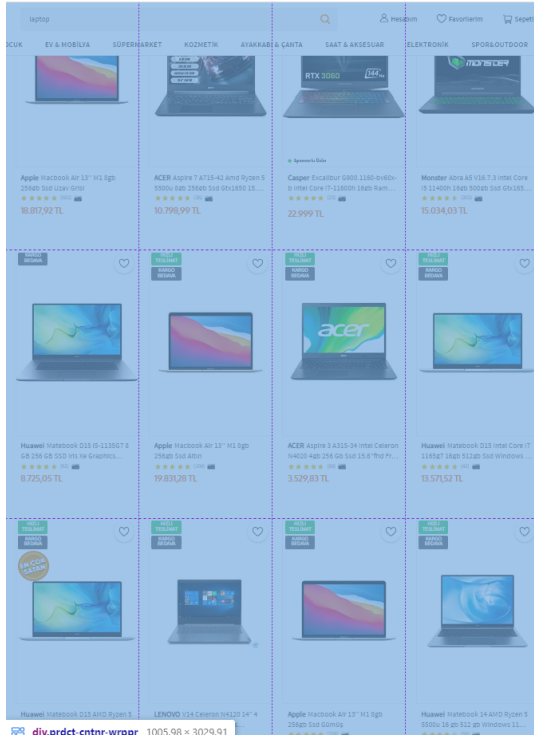
```

laptop_detail_link_rules = LinkExtractor(restrict_xpaths='//div[@class="prdct-cntnr-wr
ppr"]')
laptop_detail = Rule(laptop_detail_link_rules,
                      callback='parse_item',
                      follow=False)

rules = (
    laptop_detail,
)

```

`laptop_detail` adından bir local variable oluşturduktan sonra içerisine sayfa içerisinde takip etmesini istediğim linklerin xpathini verdim.



callback='parse_item'

Sayfa içerisinden bir bilgi çıkartmak istiyorsak callback'i kullanmamız gerekiyor. Varsayılan olarak parse_item değerini alıyor.

follow=False

Eğer bu değer True olursa ziyaret edip bilgi çektiği sayfaların içerisinde LinkExtractor(restrict_xpaths="//div[@class='prdt-ctnr-wrppr']") kuralına uyan linkler varsa onları da işleme dahil edecek.