

DATA MINING

ASSIGNMENT: MILESTONE 1

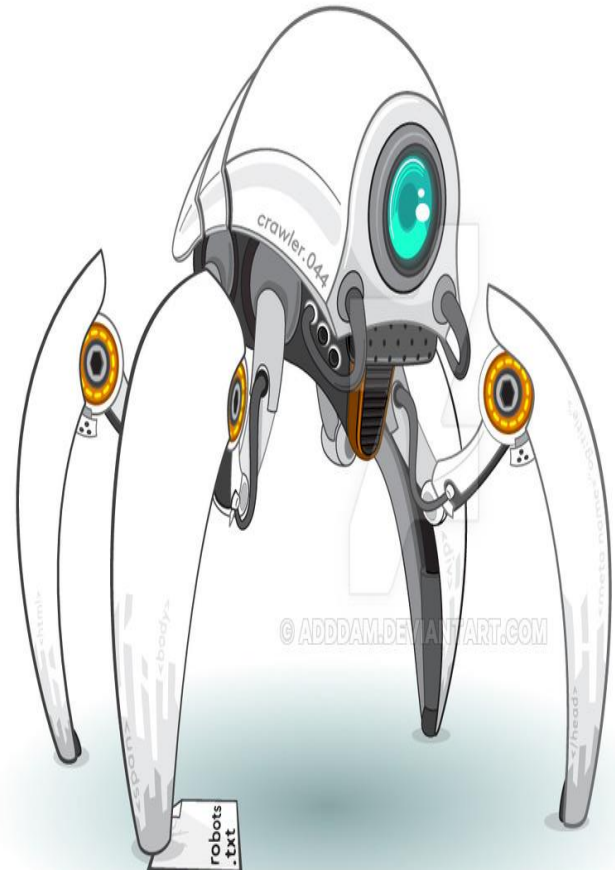
NORAISHA YUSUF (WQD180008)

NORBAIZURA MOHAMAD (WQD180043)

*MARINA SHAH MUHAMMAD ZABRI TAN
(WQD180011)*

*AMIRA AN-NUR BINTI RUSLI
(WQD180016)*

*AMINAH SOFIA BINTI MAHAYUDIN
(WQD180032)*



Stock list

source: The Star

- Main modules: Selenium and BeautifulSoup
- Nested 'for' loop
 - ✓ Outer loop: To get urls of all companies, in all markets, in alphabetical order
 - ✓ Inner loop: To crawl stock data of each specific company
- Execute_script - to extract table of list of stocks in java script

Bursa Malaysia stock quotes | T X +

https://www.thestar.com.my/business/marketwatch/stock-list/?alph

Our Sites ▾ More ▾ iBilik Propwall StarProperty.my dimsum Events Suria FM 988 FM

Star ONLINE News Business Sport Metro Tech Lifestyle Opinion Videos Property Job

ETF: Exchange Traded Fund-Bond Exchange Traded Fund-Equity Exchange Traded Fund-Commodity

View by alphabet:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z 0-9

Active	Gainers	% Gainers	Losers	% Losers			
Symbol	Open	High	Low	Last	Chg	%Chg	Vol ('00)
A50CHIN-C22	0.265	0.265	0.230	0.245	-0.045	-15.52	6,905
A50CHIN-C24	0.865	0.865	0.860	0.860	-0.050	-5.49	200
A50CHIN-C26	0.000	0.000	0.000	0.505	0.000	0.00	0
A50CHIN-C28	0.280	0.280	0.260	0.270	-0.020	-6.90	7,412

TOPICS ▸ Asean+ True or Not Do You Know Star Golden Hearts Award SOBA 2018

Stocks

Thursday, 7 March 2019

View by alphabet:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z 0-9

A50CHIN-C26: CW ISHARES FTSE A50 CHINA INDEX ETF (RHB)

Board : Warrants 52 Week High : 0.505
Stock Code : 070326 52 Week Low : 0.185

Open	High	Low	Last	Chg	Chg %	Vol ('00)	Buy/Vol ('00)	Sell/Vol ('00)
-	-	-	0.505	0.000	0.00	0	0.495 / 3,000	0.505 / 3,000

Updated : 07 Mar 2019 | 2:45 PM

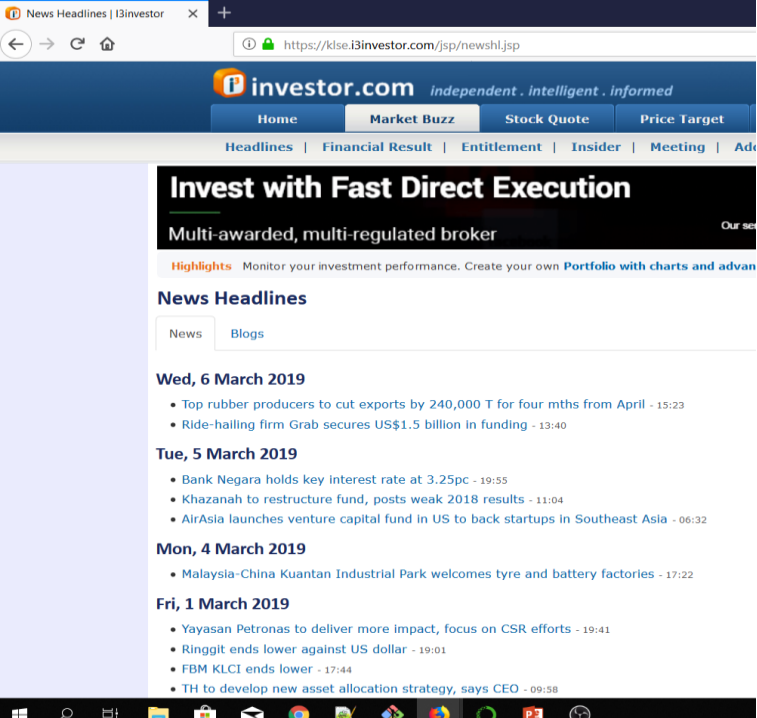
aisha_datamining_crawl.py

```
36 for i in alpha:
37     print("!!! Now char "+ i)
38     browser = webdriver.Firefox(executable_path=r'C:\geckodriver.exe')
39     browser.get(urlTheStar + i)
40     innerHTML = browser.execute_script('return document.body.innerHTML')
41     soup = BeautifulSoup(innerHTML, 'lxml')
42     stock_table = soup.find('table',{'class':'market-trans'})
43     links = stock_table.findAll('a')
44
45     company = []
46     for link in links:
47         start_page = requests.get('https://www.thestar.com.my'+link.get('href'))
48         tree = html.fromstring(start_page.text)
49
50         url_link = 'https://www.thestar.com.my'+link.get('href')
51         board = tree.xpath('//li[@class="f14"]/text()')[0]
52         stock_code = tree.xpath('//li[@class="f14"]/text()')[1]
53         name = tree.xpath('//h1[@class="stock-profile f16"]/text()')[0]
54         w52high = tree.xpath('//li[@class="f14"]/text()')[2]
55         w52low = tree.xpath('//li[@class="f14"]/text()')[3]
56         updateDate = tree.xpath('//span[@id="slcontent_0_ileft_0_datetxt"]/text()')[0]
57         updateTime = tree.xpath('//span[@class="time"]/text()')[0]
58         open_price = tree.xpath('//td[@id="slcontent_0_ileft_0_lastdonetxt"]/text()')[0]
59         high_price = tree.xpath('//td[@id="slcontent_0_ileft_0_opentext"]/text()')[0]
60         low_price = tree.xpath('//td[@id="slcontent_0_ileft_0_lowtext"]/text()')[0]
61         last_price = tree.xpath('//td[@id="slcontent_0_ileft_0_lastdonetxt"]/text()')[0]
62         volume = tree.xpath('//*[@id="slcontent_0_ileft_0_voltext"]/text()')[0]
63         buy_vol_hundred = tree.xpath('//*[@id="slcontent_0_ileft_0_buyvol"]/text()')[0]
64         sell_vol_hundred = tree.xpath('//*[@id="slcontent_0_ileft_0_sellvol"]/text()')[0]
65         date_crawl = str(datetime.datetime.now())
66
```

aisha_datamining_crawl.py

```
0.360
0.085
AHMAD ZAKI RESOURCES BERHAD- WA 14/24
07 Mar 2019
```

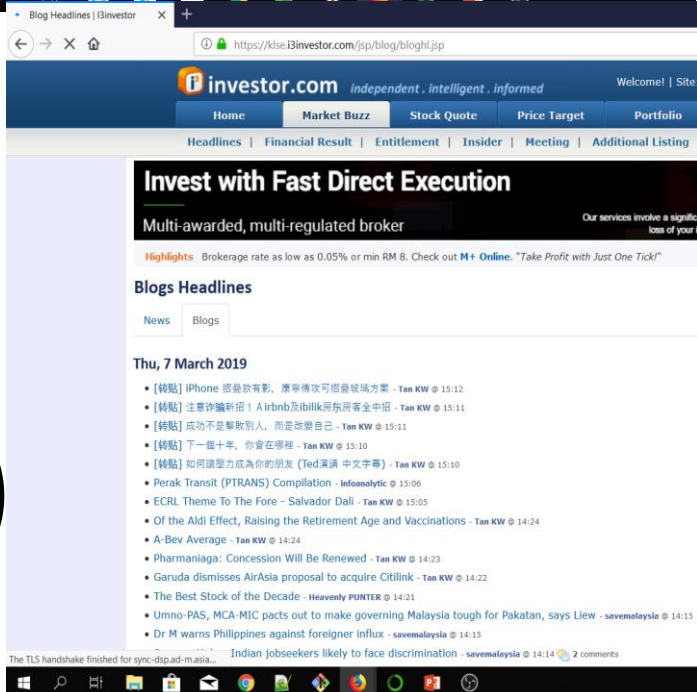
[]:



News & Blogs

source: *Investor.com*

- Selenium and BeautifulSoup
- Nested for loops
- Date, headlines and author



JupyterLab

localhost:8892/lab

File Edit View Run Kernel Tabs Settings Help

aisha_datamining_crawl.py

```
96
97 #####
98 # CRAWLING NEWS
99 #####
100 print("!!! START CRAWLING NEWS")
101 urlNews='https://klse.i3investor.com/jsp/newshl.jsp'
102 browser = webdriver.Firefox(executable_path=r'C:\geckodriver.exe')
103 browser.get(urlNews)
104 innerHTML = browser.execute_script('return document.body.innerHTML')
105 soup = BeautifulSoup(innerHTML, 'lxml')
106 date = soup.select('div > h3')
107 for a in date:
108     print(" ")
109     print(a.text)
110     div = soup.find('h3', text=a.text).find_next_siblings('ul')[0]
111     title = div.find_all('a')
112     for b in title:
113         time_raw = b.find_next_siblings('span', {'class': 'graydate'})[0].text
114         time = time_raw[3:].strip()
115         dateInsert = str(datetime.datetime.now())
116         tarikh = a.text
117         tajuk = b.text
118         penulis = None
119         category = "news"
120         print(b.text)
121         print(time)
```

aisha_datamining_crawl.py

```
07:55

Tue, 5 February 2019
Germania airline says filed for bankruptcy, cancels all flights
23:43
Merck KGaA wins GSK for immunotherapy deal worth up to $4.2 bln
20:48
Volkswagen courting Swedish investors to anchor Traton truck IPO
19:50
```

[]:

Windows Taskbar

8:36 PM 7/3/2019

JupyterLab

localhost:8892/lab

File Edit View Run Kernel Tabs Settings Help

aisha_datamining_crawl.py

```
138 #soup = BeautifulSoup(page.text, 'html.parser')
139 innerHTML = browser.execute_script('return document.body.innerHTML')
140 soup = BeautifulSoup(innerHTML, 'lxml')
141 date = soup.find("div", {"id": "maincontent730"}).find_all('h3')
142 print(date)
143 for a in date:
144     print(" ")
145     data_ul = soup.find('h3', text=a.text).find_next_siblings('ul')[0]
146     #data_li = data_ul.select('ul > li')
147     data_li = data_ul.findAll('li')
148     for b in data_li:
149         title = b.find('a')
150         author = b.find('span', {'class': 'comuid'})
151         all_text = b.find('span', {'class': 'graydate'}).text
152         child_text = b.find('span', {'class': 'comuid'}).text
153         parent_text = all_text.replace(child_text, '')
154         print(" ")
155         dateInsert = str(datetime.datetime.now())
156         tarikh = a.text
157         print(tarikh)
158         tajuk = title.text
159         print(tajuk)
160         penulis = author.text
161         print(penulis)
162         category = "blog"
163         time = parent_text[5:].strip()
```

aisha_datamining_crawl.py

```
Thu, 28 February 2019
YTL Hosp REIT - New Assets Cushioned Weaker AUD Earnings
kltrader

Thu, 28 February 2019
Evening Market Summary - 27 Feb 2019
mplus313
```

[]:

Windows Taskbar

8:37 PM 7/3/2019

Financial Performance

source: Investor.com

- Mainly Selenium module
- For Loop within While statement function
- Element.click() – to automate clicking of check boxes and table page index

investor.com X +

https://klse.i3investor.com/financial/quarter/latest.jsp

----- Price -----

☐ Price ☐ Change ☐ %

----- Financial Result -----

☐ Revenue ☐ PBT ☐ NP ☐ NP to SH ☐ Div

----- Financial Ratio -----

☐ Div Payout % ☐ NP Margin ☐ ROE

----- Per Share Item -----

☐ RPS ☒ EPS ☒ DPS ☐ NAPS

----- Performance -----

☒ QoQ ☒ YoY

Show 25 entries Search:

Stock	Date			Per Share Item		QoQ
	Ann. Date	F.Y.	Quarter	EPS	DPS	
ETH	05-Mar-2019	31-Dec-2018	31-Dec-2018	-0.53	0.00	- %
DESTINI	05-Mar-2019	31-Dec-2018	31-Dec-2018	0.14	0.00	↑ 111.90%
MUHIBAH	05-Mar-2019	31-Dec-2018	31-Dec-2018	7.87	7.50	↑ 0.37%
UOADEV	05-Mar-2019	31-Dec-2018	31-Dec-2018	7.40	14.00	↑ 48.05%
T7GLOBAL	04-Mar-2019	31-Dec-2018	31-Dec-2018	1.16	0.00	↑ 357.80%
PPB	01-Mar-2019	31-Dec-2018	31-Dec-2018	15.56	20.00	↓ -38.4%
MBSB	01-Mar-2019	31-Dec-2018	31-Dec-2018	1.88	0.00	↓ -3.0%
SANICHI	01-Mar-2019	31-Dec-2018	31-Dec-2018	-2.76	0.00	↓ -54.0%
KPOWER	01-Mar-2019	30-Jun-2019	31-Dec-2018	-3.89	0.00	↓ -243.90%

JupyterLab

localhost:8892/lab

FileEditViewRunKernelTabsSettingsHelp

aisha_datamining_crawl.py

178# CRAWLING Financial Info

179#####

180print("!!! START CRAWLING FINANCIAL DATA")

181

182browser = webdriver.Firefox(executable_path=r'C:\geckodriver.exe')

183browser.implicitly_wait(40)

184

185Financialurl = 'https://klse.i3investor.com/financial/quarter/latest.jsp'

186browser.get(Financialurl)

187

188#to expand the "modify the visible columns i.e. checkboxes"

189WebElementexpanded = browser.find_element_by_xpath("//*[@id='ui-accordion-financialResultTableColumnsDiv-header-0']/span")

190WebElementexpanded.click()

191

192# to ensure all checkboxes are checked

193alllinks = browser.find_elements_by_xpath('//input[@type="checkbox"]')

194for link in alllinks:

195 if link.is_selected():

196 print('Checkbox already selected');

197 else:

198 link.click();

199 print('Checkbox selected');

aisha_datamining_crawl.py

Checkbox selected

Checkbox selected

Checkbox selected

Checkbox already selected

Checkbox already selected

Checkbox selected

Checkbox already selected

Checkbox already selected

 aisha_datamining_crawl.py

[]: