

Webscraping with Selenium and BeautifulSoup

Thomas Jonas

2019-06-16

Ziel

- ▶ Automatisiertes Webscraping
- ▶ Datenhaltung in Zeitreihen-Datenbank
- ▶ Darstellung durch Dashboard
- ▶ Deployment in Docker-Umgebung

Inhaltsangabe

- ▶ Übersicht genutzter Technologien
- ▶ Implementierung & Konfiguration
- ▶ Probleme
- ▶ Fazit

Webscraping

Selenium

- ▶ Webscraping von dynamischen single-page Webseiten.
- ▶ Nutzt Browser (chrome-driver muss installiert sein.)

BeautifulSoup

- ▶ Suche nach Klassen und Tags innerhalb HTML-Code.
- ▶ Light-weight

Zeitreihen Datenbank

influxDB

- ▶ Schreiben/Lesen von Zeitreihen-Daten.
- ▶ Meist genutzte ZR-DB.
- ▶ Push-basiert.

Prometheus

- ▶ Sammelt Metriken über zB CPU-Auslastung
- ▶ Pull-basierte Datenbank

Dashboard

Grafana

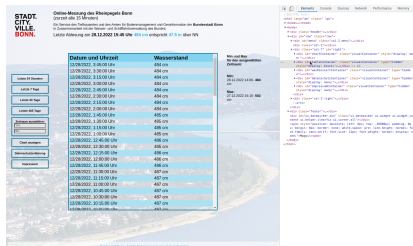
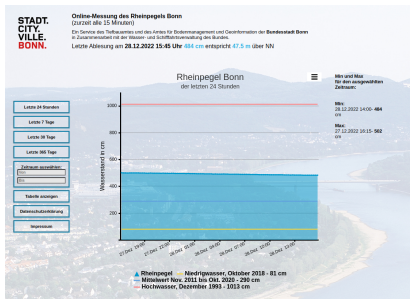
- ▶ Interactive Web-application.
- ▶ Erstellen anschaulicher Graphen und Diagramme.
- ▶ Mehrere Datenquellen möglich.

Deployment

Docker

- ▶ OS-level Virtualisierung. Verpacken von Software in Container.

Implementierung: Webscraping



Implementierung: Webscraping

[illegible]

- ▶ Selenium ermöglicht Navigation auf Webseiten.
- ▶ Auswahl, Ausführung und Änderung von Elementen.

Implementierung: Webscraping - Selenium

```
1 service_object = Service(binary_path)
2 # Invoke new browser window & enable headless mode
3 chrome_options = Options()
4 chrome_options.add_argument("--headless")
5 chrome_options.add_argument("--window-size=1920x1080")
6 chrome_options.add_argument('--no-sandbox')
7 chrome_options.add_argument('--disable-gpu')
8 driver = webdriver.Chrome(service=service_object, options=chrome_options)
9
10 # Navigate to website
11 driver.get("https://pegel.bonn.de/php/rheinpegel.php")
12 # Select button element
13 show_table_button = driver.find_element(by='id', value='btn_table')
14 # Click button
15 show_table_button.click()
16 # Retrieve generated HTML element by id
17 waterlevel_data_element = driver.find_elements(by='id', value='dataTable')
18 # Get innerHTML content of table
19 waterlevel_data_html = waterlevel_data_element[0].get_attribute(
20     name='innerHTML')
21 driver.close()
```

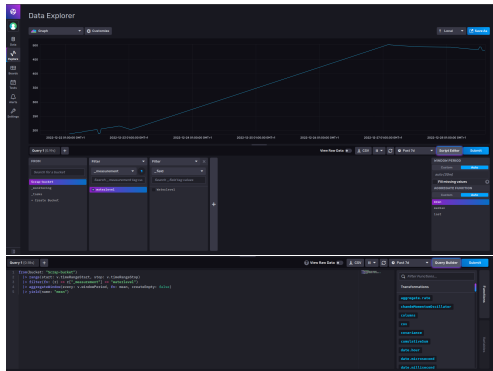
Implementierung: Webscraping - BeautifulSoup

```
1  # BeautifulSoup used to retrieve elements inside specific HTML-tags.
2
3  # Create Soup from HTML
4  soup = BeautifulSoup(waterlevel_data_html, 'html.parser')
5  # Find all table rows
6  table_rows = soup.find_all('tr')
7  table_header = table_rows[0].text
8  table_data = table_rows[1:]
9  # Split into single td (tabledata) elements
10 html_rows = [BeautifulSoup(str(table_data[i]), 'html.parser').find_all('td')
11               for i in range(len(table_data))]
12 rows = [[row[0].text, row[1].text] for row in html_rows]
13 date_time_raw = [row[0] for row in rows]
14 waterlevel_raw = [row[1] for row in rows]
15
16 return date_time_raw, waterlevel_raw
```

Implementierung: Datenhaltung - influxDB

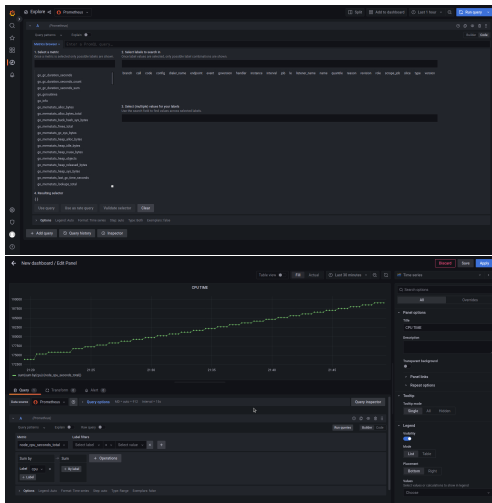
```
1 def map_duplicates_to_24_hour_format(dataFrame):
2     ...
3
4 def set_index_column_time(dataframe):
5     ...
6
7 def write_to_influxdb(dataframe):
8     client = InfluxDBClient(url=f"http://{INFLUXDB_URL}:8086",
9                             token=influx_token, org=influx_org)
10    write_api = client.write_api(write_options=SYNCHRONOUS)
11    write_api.write(influx_bucket, influx_org, record=dataframe,
12                   data_frame_measurement_name='waterlevel',
13                   data_frame_tag_columns=['waterlevel'])
14    write_api.close()
15    client.close()
```

Dashboard - influxDB



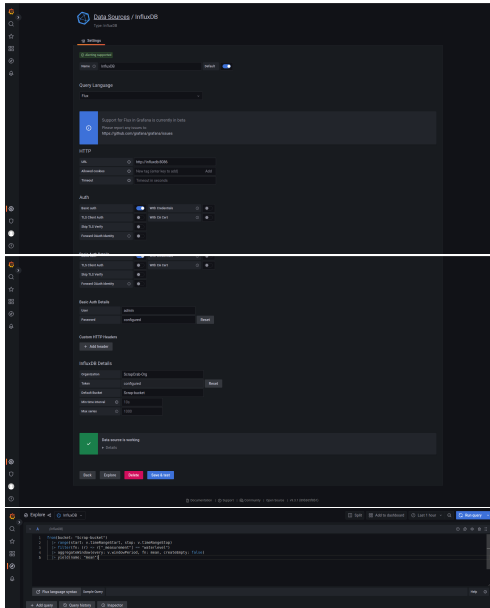
- ▶ Query Builder
- ▶ Script Editor (hilfreich für Grafana)

Dashboard - Grafana



- ▶ Prometheus als Datenquelle hinzugefügt.
- ▶ Monitoring unterschiedlicher Metriken. zB CPU-Auslastung

Dashboard - Grafana



- ▶ influxDB als Datenquelle hinzufügen.
- ▶ Credentials durch .env-Variablen in influxDB container gesetzt.
- ▶ Query aus influxDB-UI einfügen.

Dashboard - Grafana



Deployment: docker-compose

```
1  version: '3.8'
2  services:
3    influxdb:
4      container_name: influxdb
5      image: influxdb:2.0
6      volumes:
7        - ./influx/database:/var/lib/influxdb2
8      env_file:
9        - .env
10     environment:
11       - DOCKER_INFLUXDB_INIT_MODE=setup
12       - DOCKER_INFLUXDB_INIT_USERNAME=${influx_user}
13       - DOCKER_INFLUXDB_INIT_PASSWORD=${influx_pw}
14       - DOCKER_INFLUXDB_INIT_ORG=${influx_org}
15       - DOCKER_INFLUXDB_INIT_BUCKET=${influx_bucket}
16       - DOCKER_INFLUXDB_INIT_ADMIN_TOKEN=${influx_token}
17
18     ports:
19       - "8086:8086"
20
21     web-scraper:
22       container_name: scrapcrab
23       image: scrapcrab:0.1
24       env_file:
25         - .env
26       environment:
27         - INFLUXDB_URL=influxdb
28         - REFRESH_INTERVAL_MINUTES=30
29       command: python scheduled_scraper.py
30       depends_on:
31         - influxdb
```

```
1  prometheus:
2    image: prom/prometheus:latest
3    container_name: prometheus
4    restart: unless-stopped
5    volumes:
6      - ./prometheus.yml:/etc/prometheus/prometheus.yml
7      - prometheus_data:/prometheus
8    command:
9      - '--config.file=/etc/prometheus/prometheus.yml'
10     - '--storage.tsdb.path=/prometheus'
11     - '--web.console.libraries=/etc/prometheus/console_libraries'
12     - '--web.console.templates=/etc/prometheus/consoles'
13     - '--web.enable-lifecycle'
14
15    ports:
16      - 9090:9090
17
18    grafana:
19      image: grafana/grafana-oss:9.3.1-ubuntu
20      container_name: grafana
21      ports:
22        - 3000:3000
23      volumes:
24        - ./grafana:/var/lib/grafana
25      user: "0:0"
26      environment:
27        GF_SECURITY_ADMIN_PASSWORD: grafana_totaly_secure_password
28
29    volumes:
30      prometheus_data: {}
```

Probleme

Docker-Images Zugriffsrechte

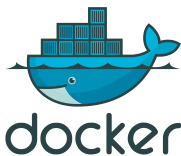
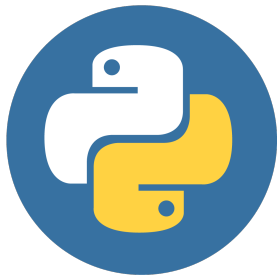
```
prometheus:
  image: prom/prometheus:latest
  container_name: prometheus
  restart: unless-stopped
  volumes:
    - ./prometheus.yml:/etc/prometheus/prometheus.yml
    - ./prometheus_data:/prometheus
  command:
```

```
prometheus | ts=2022-12-20T14:00:02.595Z caller=query_logger.go:91 level=error component=activeQueryTracker msg="Error opening query log file" file=/
prometheus | panic: open /etc/prometheus/queries/active: permission denied
prometheus | panic: Unable to create mmap-ed active query log
prometheus |
prometheus | goroutine 1 [running]:
prometheus | github.com/prometheus/prometheus/promql.NewActiveQueryTracker(0x77fcb2bdef7, 0xb0, 0x14, 0x30a20eb, 0xc00180d093)
prometheus | /app/promql/query_logger.go:321 +0x3cd
prometheus | main.main()
prometheus | /app/cmd/prometheus/main.go:618 +0x6973
prometheus |
prometheus exited with code 2
```

```
- prometheus_data:/prometheus
```

- ▶ Fehlende Zugriffsrechte verursache bei Prometheus/ Grafana Images oft Probleme.
- ▶ Lösung: Docker Volumes

Genutzte Technologien



BeautifulSoup



TODO

- ▶ Page-Nummern einfuegen
- ▶ Datum und Author auf jeder Folie