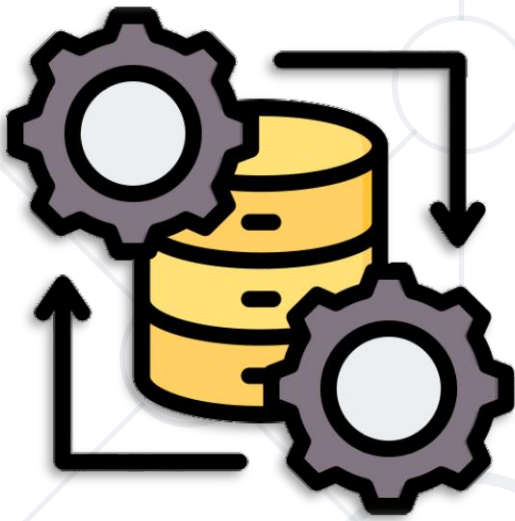


Data Science in Production

Data Pipelines and Processes



Yordan Darakchiev
Technical Trainer



SoftUni



Software University

<https://softuni.bg>

Have a Question?

sli.do

#DataScience

1. {Dev, Data, ML}Ops Fundamentals
2. Work Processes
3. Data Pipelines
4. Scaling Workflows to Larger Datasets
5. Monitoring and Dashboards





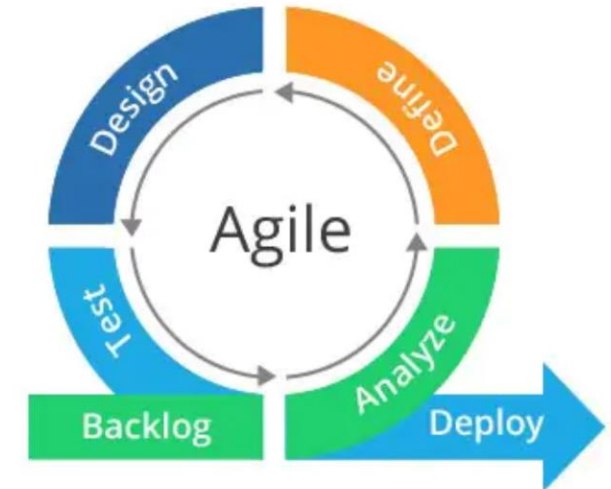
Work Processes

Collaborating with Others

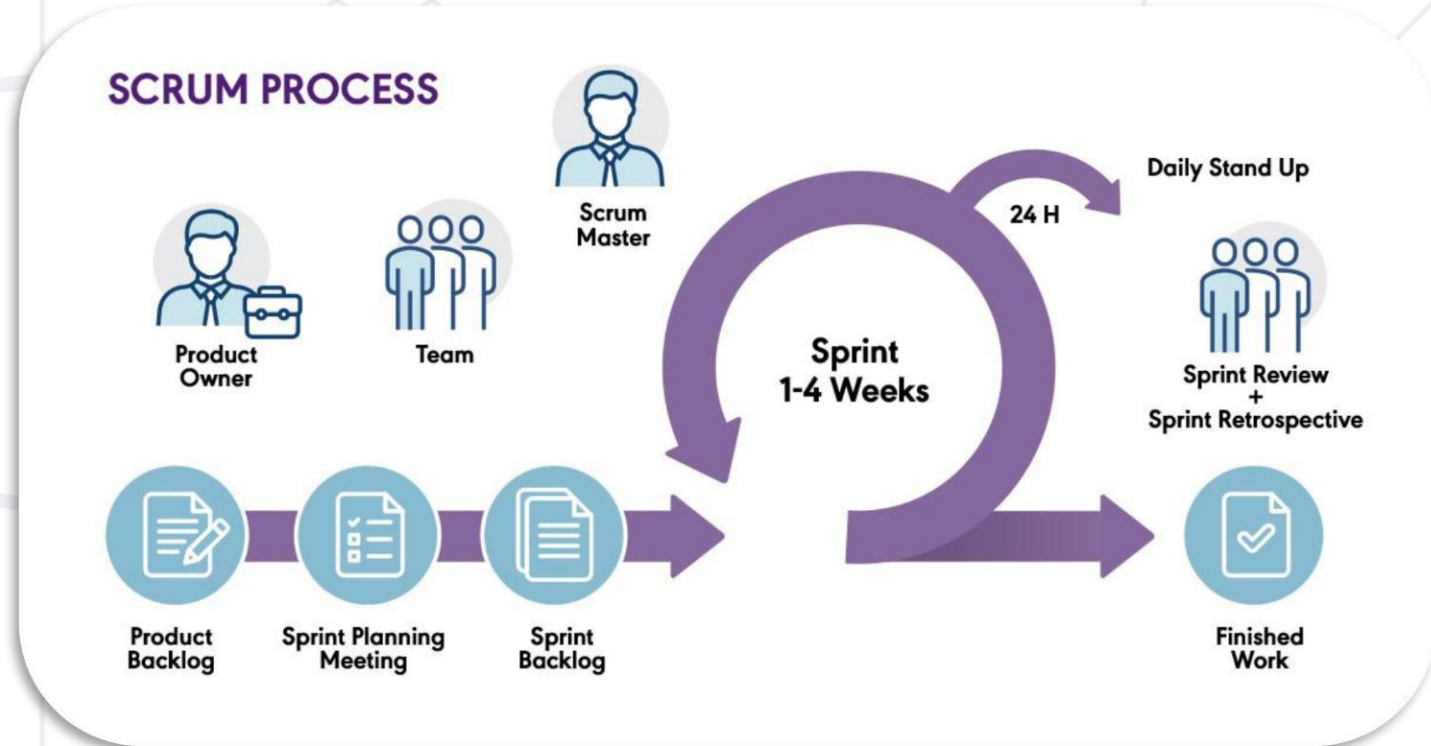
Agile (vs. Waterfall)

- Incremental
- **Iterative**
- Focused on what our client needs
- *The Agile Manifesto*

Waterfall vs. Agile



- Lightweight
- Experimental
 - Figure out as you go
- **Small pieces of work at a time**
- Feedback loop
- It's not hard to get started
- Scrum Master ↔ Product Owner ↔ Scrum team
- Product backlog ↔ sprint



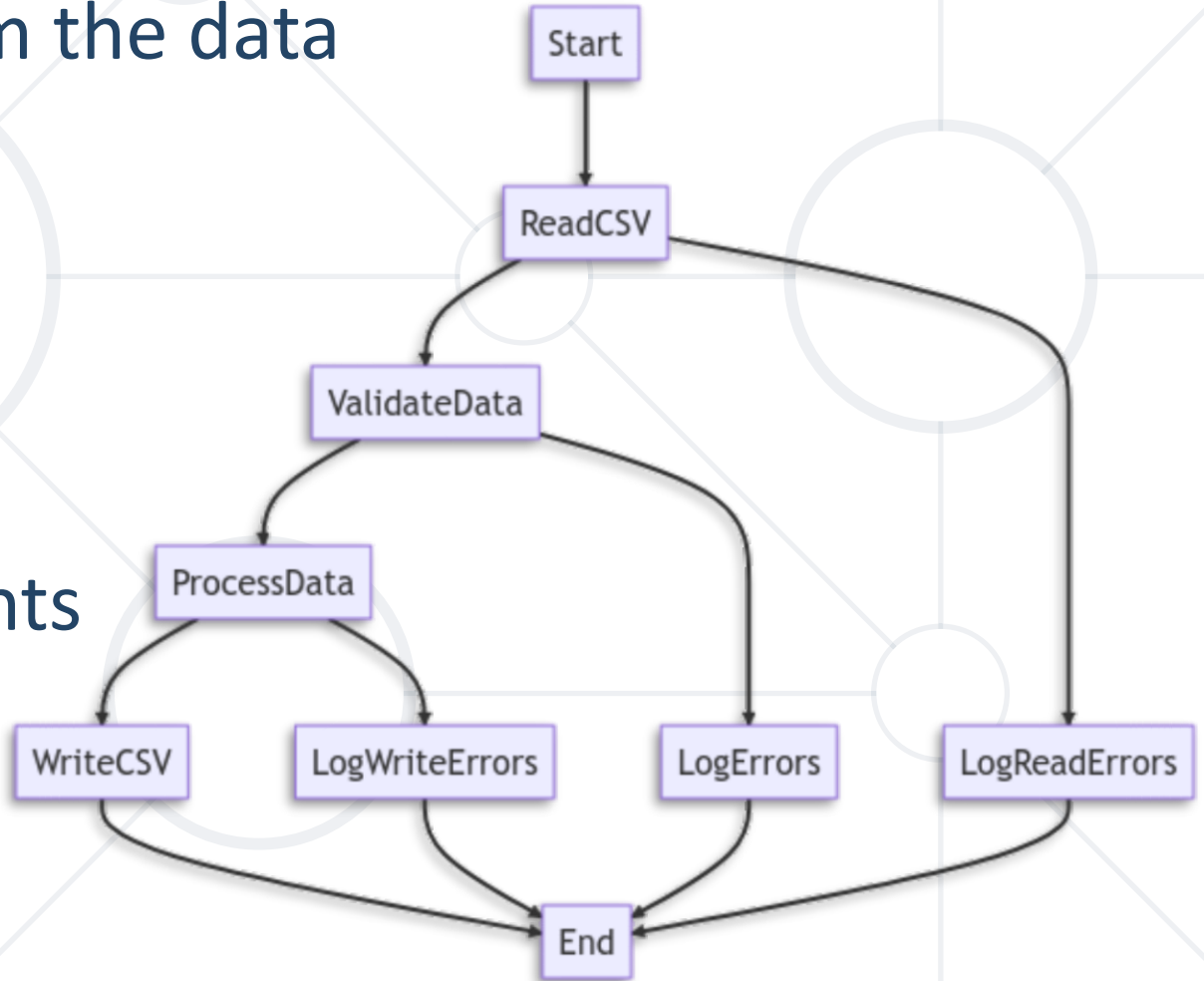


Data Pipelines

From Notebooks to Automation

Data Pipelines

- DAG of functions which transform the data
 - Automated
 - Reusable
 - Testable / tested
- Can be arbitrarily complex
- Can be used to quickly gain insights
- Can run anywhere
 - Cloud pipelines



- Installation

```
conda install anaconda::luigi
```

- Tasks and targets

- Getting started

```
class HelloLuigi(luigi.Task):  
    def output(self):  
        return luigi.LocalTarget('hello-  
luigi.txt')  
  
    def run(self):  
        with self.output().open("w") as outfile:  
            outfile.write("Hello Luigi!")
```

```
python -m luigi --module <file name> HelloLuigi
```

■ Parameters

```
class ParamDemo(luigi.Task):  
    my_param = luigi.IntParameter(  
        42,  
        visibility=luigi.parameter.ParameterVisibility.PRIVATE)  
  
    def output(self):  
        print(self.my_param) # Usage: directly, with self
```

■ Configuration

- luigi.cfg
- Parameter values can be overwritten (code config overwrites file config)



Big(ger) Data

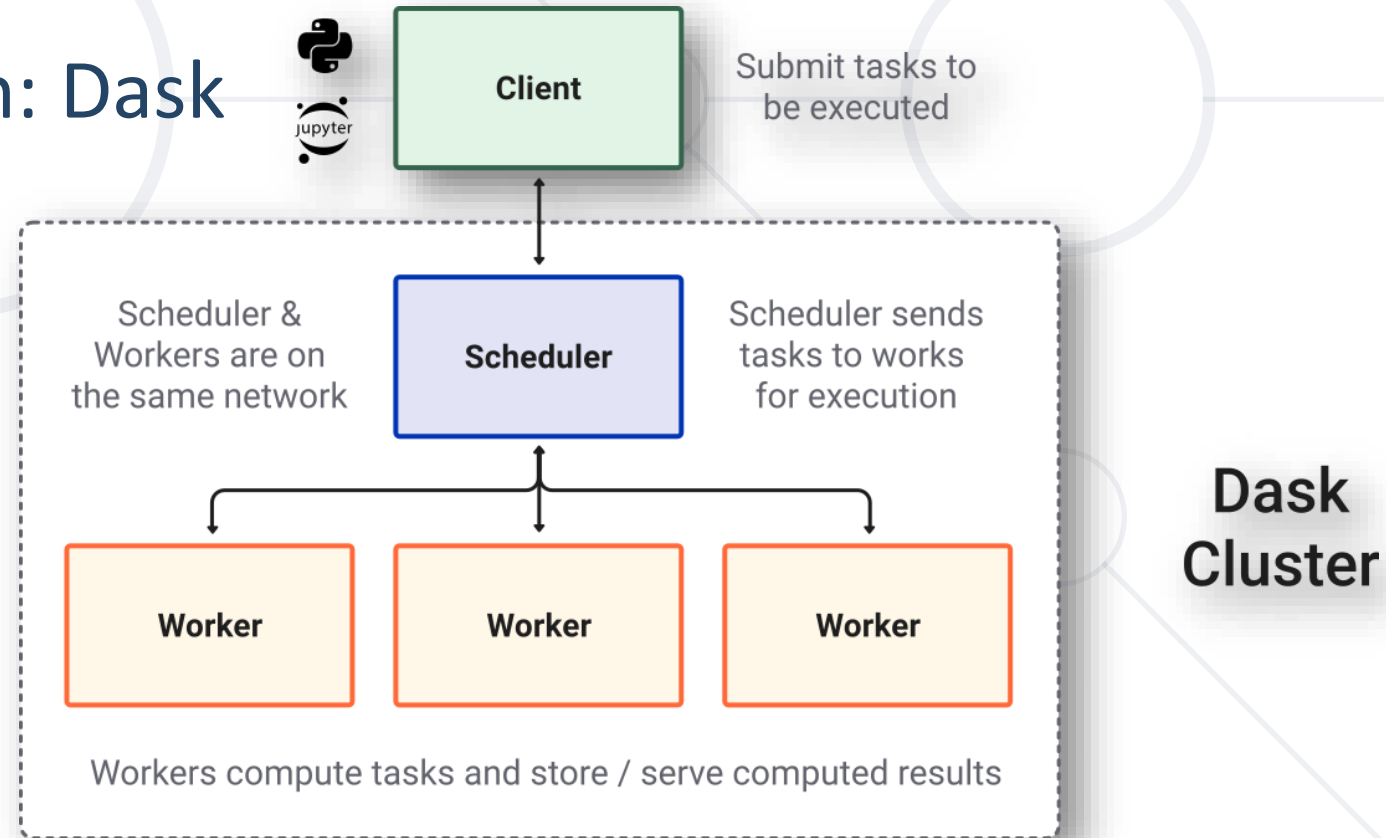
Scaling up our Scripts and Computation

Getting to Production

- Data versioning: DVC
- Experiment tracking: MLFlow
- Concurrent data manipulation: Dask

```
import dask

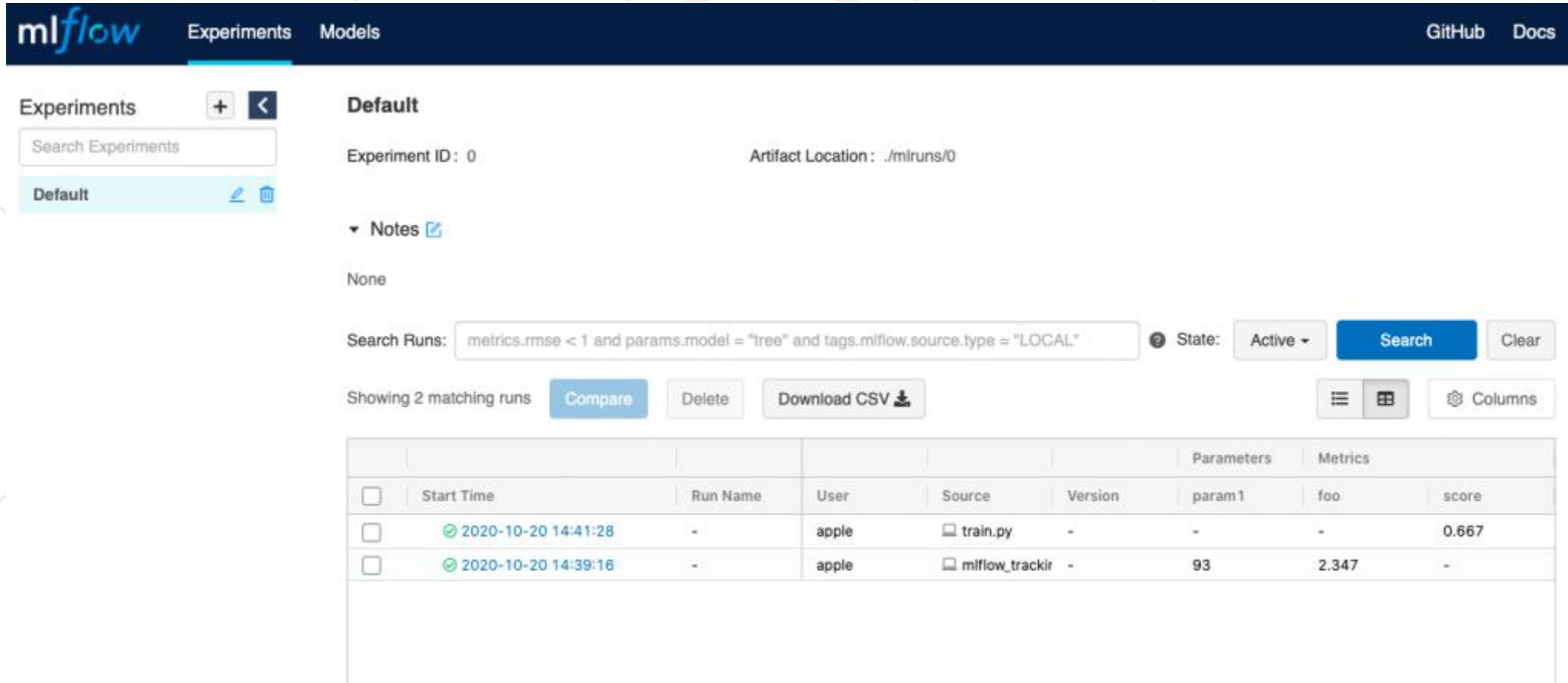
df = dask.datasets.timeseries()
print(df[df.y > 0])
```





Monitoring
Getting Insights

- MLFlow tracking can do model / data / artifact monitoring



The screenshot displays the MLFlow web interface for monitoring experiments. The top navigation bar includes the MLFlow logo, tabs for 'Experiments' and 'Models', and links to 'GitHub' and 'Docs'. On the left, a sidebar shows the 'Experiments' section with a search bar and a list of experiments, including 'Default'. The main content area is titled 'Default' and shows 'Experiment ID: 0' and 'Artifact Location: ./mlruns/0'. Below this, there is a 'Notes' section which is currently empty. A search bar is present with the query 'metrics.rmse < 1 and params.model = "tree" and tags.mlflow.source.type = "LOCAL"'. The search results show 2 matching runs. Below the search bar, there are buttons for 'Compare', 'Delete', and 'Download CSV'. A table displays the search results with columns for Start Time, Run Name, User, Source, Version, Parameters, and Metrics.

| | Start Time | Run Name | User | Source | Version | Parameters | Metrics |
|--------------------------|---------------------|----------|-------|----------------|---------|------------|---------|
| <input type="checkbox"/> | 2020-10-20 14:41:28 | - | apple | train.py | - | - | 0.667 |
| <input type="checkbox"/> | 2020-10-20 14:39:16 | - | apple | mlflow_trackir | - | 93 | 2.347 |

- Many libraries can do that
- Plotly's dash is really popular
 - Examples, tutorial

```
from dash import Dash, html, dash_table, dcc
import pandas as pd
import plotly.express as px

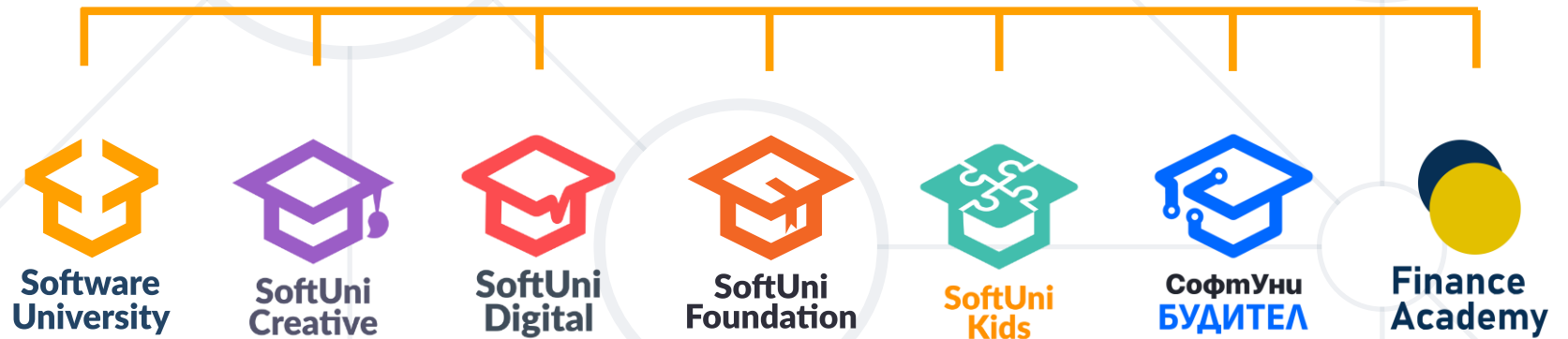
df = pd.read_csv("...")
app = Dash()
app.layout = [
    html.Div(children="My First App with Data and a Graph"),
    dash_table.DataTable(data=df.to_dict("records"), page_size = 10),
    dcc.Graph(figure = px.histogram(df, x = "continent", y = "lifeExp",
histfunc = "avg"))
]
if __name__ == "__main__":
    app.run(debug = True)
```

Summary

- {Dev, Data, ML}Ops Fundamentals
- Work Processes
- Data Pipelines
- Scaling Workflows to Larger Datasets
- Monitoring and Dashboards



Questions?



SoftUni Diamond Partners



- Software University – High-Quality Education, Profession and Job for Software Developers
 - softuni.bg, about.softuni.bg
- Software University Foundation
 - softuni.foundation
- Software University @ Facebook
 - facebook.com/SoftwareUniversity



- This course (slides, examples, demos, exercises, homework, documents, videos and other assets) is **copyrighted content**
- Unauthorized copy, reproduction or use is illegal
- © SoftUni – <https://about.softuni.bg/>
- © Software University – <https://softuni.bg>

