# Productionizing Jupyter Notebooks

Antoni Ivanov

Versatile Data Kit Open Source Project

VERSATILE DATA KIT

# The role of Jupyter in the data world



https://analyticsindiamag.com/why-jupyter-notebooks-are-so-popular-among-data-scientists/

https://odsc.medium.com/why-you-should-be-using-jupyter-notebooks-ea2e568c59f2

# The role of Versatile Data Kit (VDK) in the data world

Develop

Deploy and Monitor



VDK SDK

```
extract_load_rest_calls.py

def run(job_input):
    response = requests.get("https://rest.com/calls")
    payload = response.json()

    job_input.send_object_for_ingestion(
        payload=payload,
        destination_table="rest_target_table")
```

```
transform_sales_mart.sql

insert into {mart_schema}.{sales_table}
SELECT
    s.product_id,
    s.transaction_date,
    s.quantity_sold * p.product_price
FROM {raw_schema}.{sale_transaction_table} as s
JOIN {raw_schema}.{products_table} p using product_id
```

Control Plane and Operations UI

https://github.com/vmware/versatile-data-kit

VERSATILE
DATA KIT

3

# Jupyter
## From Data Exploration to Production

# Challenges

➢ Reproducibility: Non-Linear Execution and Hidden State Risks

➢ Code Organization: Irrelevant or debugging code

➢ Execution model: interactive kernel vs automated flow

➢ Automated Testing and CICD

➢ Version Control

VERSATILE
DATA KIT

# Reproducibility: Non-Linear Execution and Hidden State Risks

```
co = 0
```

```
co += 1
```

```
co
```

VERSATILE
DATA KIT

# Reproducibility: Non-Linear Execution and Hidden State Risks

VERSATILE
DATA KIT

# What can we do?

```python
[2]: import pandas as pd
     # Read the data
     url = "https://raw.githubusercontent.com/duyguHsnHsn/nps-data/main/nps_data.csv"
     df = pd.read_csv(url)
```

```python
[3]: df = df[df['User'] != 'testuser']
```

```python
[4]: df.head()
```

[4]:

|   | Date | User | Score |
|---|------|------|-------|
| 1 | 2023-01-01 | mike897 | 5 |
| 2 | 2023-01-01 | lucy131 | 7 |
| 3 | 2023-01-01 | david479 | 5 |
| 4 | 2023-01-01 | david220 | 0 |
| 6 | 2023-01-02 | alex467 | 9 |

```python
[5]: job_input.send_tabular_data_for_ingestion(
         df.itertuples(index=False),
         destination_table="nps_data",
         column_names=df.columns.tolist()
     )
```

```python
[6]: %%vdksql
     select * from nps_data
```

[6]:

|   | Date | User | Score |
|---|------|------|-------|
| 0 | 2023-01-01 | mike897 | 5 |
| 1 | 2023-01-01 | lucy131 | 7 |

VERSATILE
DATA KIT

# What can we do?

## Tagging VDK Cells

```
[2]:  import pandas as pd                                                                    1
      # Read the data
      url = "https://raw.githubusercontent.com/duyguHsnHsn/nps-data/main/nps_data.csv"
      df = pd.read_csv(url)
```

```
[3]:  df = df[df['User'] != 'testuser']                                                      2
```

```
[4]:  df.head()
```

[4]:

|   | Date | User | Score |
|---|------|------|-------|
| 1 | 2023-01-01 | mike897 | 5 |
| 2 | 2023-01-01 | lucy131 | 7 |
| 3 | 2023-01-01 | david479 | 5 |
| 4 | 2023-01-01 | david220 | 0 |
| 6 | 2023-01-02 | alex467 | 9 |

```
[5]:  job_input.send_tabular_data_for_ingestion(                                             3
          df.itertuples(index=False),
          destination_table="nps_data",
          column_names=df.columns.tolist()
      )
```

```
[6]:  %%vdksql
      select * from nps_data
```

[6]:

|   | Date | User | Score |
|---|------|------|-------|
| 0 | 2023-01-01 | mike897 | 5 |
| 1 | 2023-01-01 | lucy131 | 7 |

VERSATILE
DATA KIT

# What can we do?

## Tagging VDK Cells

VERSATILE
DATA KIT

# Reproducibility: Non-Linear Execution and Hidden State Risks

- Assign a "vdk" tag and a specific number to a cell.
- The number dictates the order in which the cell will
  be executed in production.

Benefits:

- Ensures only the tagged cells are executed, and in the determined sequence.
- Clearly defining the execution order.
- Detect when the current state is diverging from expected order.
- Test easily end-to-end before deployment (as we will see)

VERSATILE
DATA KIT

# Challenges

✓ Reproducibility: Non-Linear Execution and Hidden State Risks

➢ **Code Organization: Irrelevant or debugging code**

➢ Execution model: interactive kernel vs automated flow

➢ Automated Testing and CICD

➢ Version Control

VERSATILE
DATA KIT

# Code Organization: Irrelevant or debugging code

```
[ ]: import pandas as pd
```

```
[ ]: url = "some-url"
     df = pd.read_csv(url)
```

Relevant

Irrelevant

```
[ ]: visualise(df)
```

VERSATILE
DATA KIT

# Code Organization: Irrelevant or debugging code

VDK tags to the rescue again

```
[ ]:   # Import all functions from the 'helper' module,           1
       # which contains the necessary logic for classification and data visualization
       from helper import visualize_data, classify_score
```

```
[ ]:   # Apply the classification function to the 'Score' column to determine the 'Type'    2
       # Note: this cell might fail on its first run.
       # If it does, simply run it again, and it should work as expected.
       df.loc[:, 'Type'] = df['Score'].apply(classify_score)
```

```
[ ]:   # Check the DataFrame
       df
```

```
[ ]:   # Visualise the types of users
       visualize_data(df)
```

## 5.2 Data Ingestion

```
[ ]:   # Sending data for ingestion                                 3
       job_input.send_tabular_data_for_ingestion(
           df.itertuples(index=False),
           destination_table="nps_data",
           column_names=df.columns.tolist()
       )
```

VERSATILE
DATA KIT

# Challenges

✓ Reproducibility: Non-Linear Execution and Hidden State Risks

✓ Code Organization: Irrelevant or debugging code

➢ Execution model: interactive kernel vs automated flow

➢ Automated Testing and CICD

➢ Version Control

VERSATILE
DATA KIT

# Execution model: interactive kernel vs automated flow

Bad for automation, bad for being part of a workfow

VERSATILE
DATA KIT

# Execution model: interactive kernel vs automated flow

# Execution model: interactive kernel vs automated flow

VERSATILE
DATA KIT

# Execution model: interactive kernel vs automated flow

✓ Reuse another notebook as a template (function)

```python
def run(job_input: IJobInput):
    args = dict(
        source_table="vm_new_data",
        target_table="dim_vm",
        timestamp_column="arrival_ts",
        id_column="vm_uuid",
    )
    job_input.execute_template("process-note-data-jupyter-notebook", args)
```

✓ Execute within a workflow
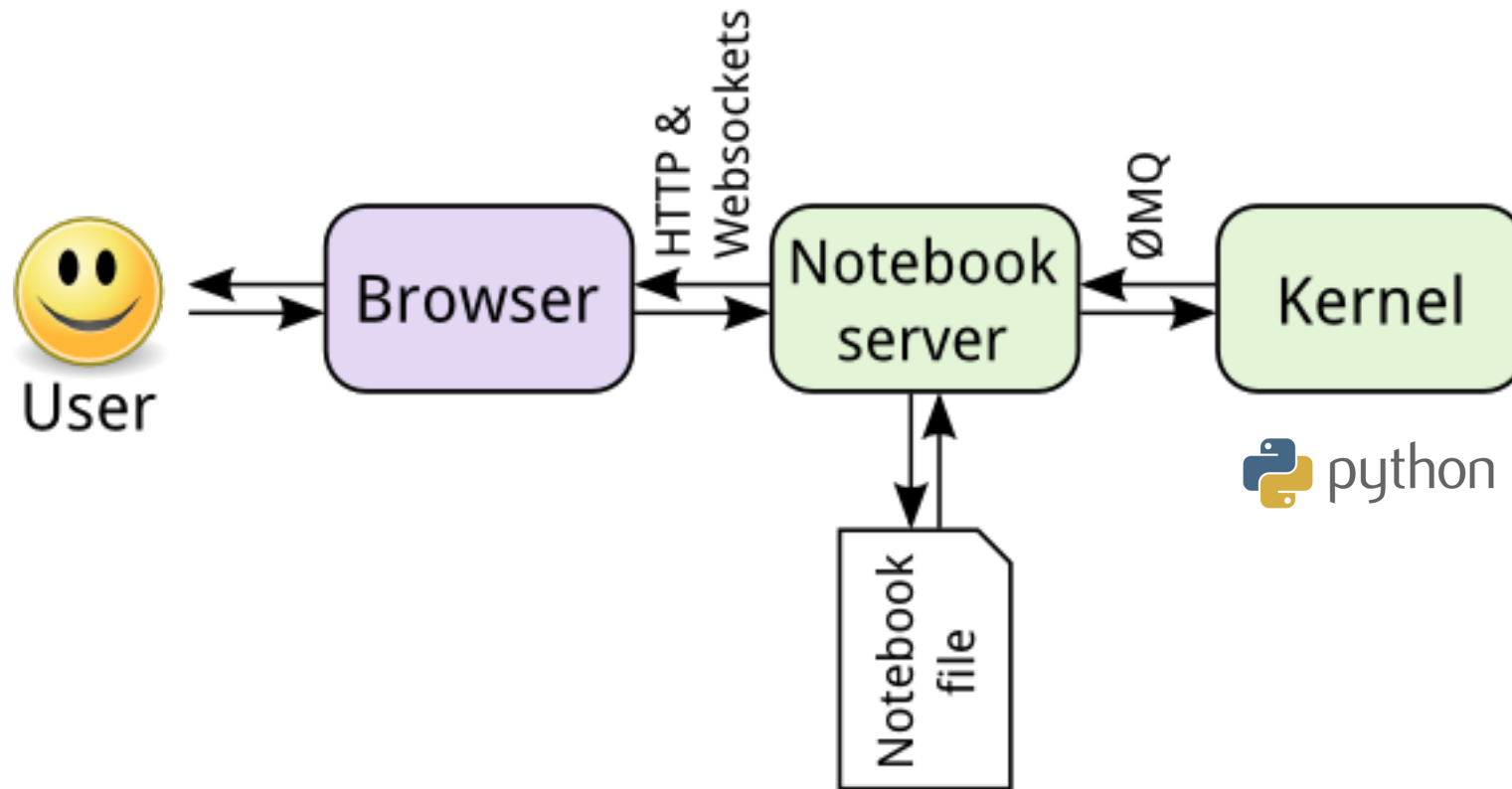
✓ Run automated tests (example coming later)

# Challenges

✓ Reproducibility: Non-Linear Execution and Hidden State Risks

✓ Code Organization: Irrelevant or debugging code

✓ Execution model: interactive kernel vs automated flow

➤ **Automated Testing and CICD**

➤ Version Control

# Automated Testing and CICD



How to Test Jupyter Notebooks with Pytest and Nbmake

Dec 14, 2021 — This tutorial describes how you can use the nbmake, a pytest plugin, to automate end-to-end **testing** of **notebooks**. jupyter **notebook** A Jupyter ...

g Notebooks Locally · Write Executed Notebooks...

Unit testing for notebooks | Databricks on AWS

How to call these functions from Python, R, Scala, and SQL **notebooks**. How to write unit **tests** in Python, R, and Scala by using the popular **test**

nteract/testbook: 🖊 📗 Unit test your Jupyter Notebooks ...

Previous attempts at unit **testing notebooks** involved writing the tests in the notebook itself. However, testbook will allow for unit tests to be run against ...

VERSATILE
DATA KIT

# Smoke (end-to-end) testing

VERSATILE
DATA KIT

# On deploy VDK requires passing smoke test first

Opt out possible.
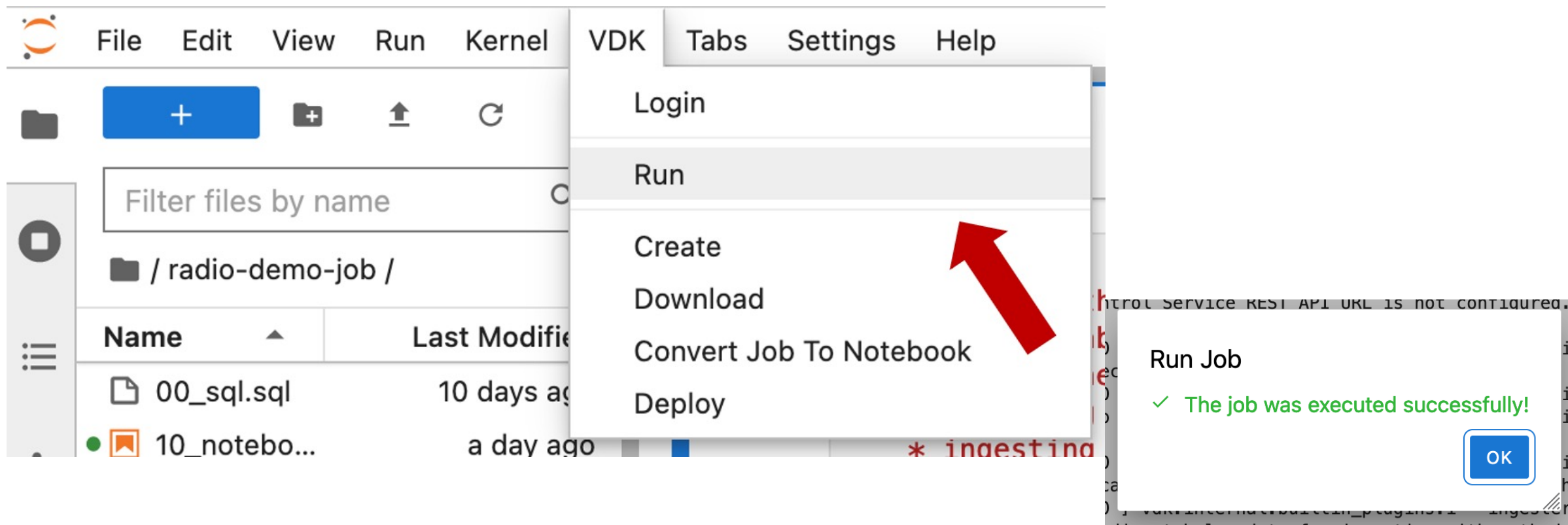
VERSATILE
DATA KIT

# Automated testing with pytest

Using VDK testing library "vdk-test-utils"

```python
from vdk.internal.test_utils import CliEntryBasedTestRunner

list_of_plugins_i_am_using = []
runner = CliEntryBasedTestRunner(list_of_plugins_i_am_using)
```

Then, invoke the data job you wish to test:

```python
result = runner.invoke(["run", "path/to/your-data-job"])
cli_assert_equal(0, result)
assert 'expected_output' in result.output
```

https://github.com/vmware/versatile-data-kit/wiki/Test-VDK-Data-Jobs-with-pytest

VERSATILE
DATA KIT

# Challenges

✓ Reproducibility: Non-Linear Execution and Hidden State Risks

✓ Code Organization: Irrelevant or debugging code

✓ Execution model: interactive kernel vs automated flow

✓ Automated Testing and CICD

➢ **Version Control**

# Version Control Challenges

```
"cell_type": "code",
"execution_count": 2,
"id": "c948f9f2-1f7b-4d8c-aeca-9b300ded9775",
"metadata": {
 "pycharm": {
  "name": "#%%\n"
 },
 "tags": [
  "vdk"
 ]
},
"outputs": [
 {
  "ename": "NameError",
  "evalue": "name 'job_input' is not defined",
  "output_type": "error",
  "traceback": [
   "\u001B[0;31m---------------------------------------------------------------------------\u001B[0m",
   "\u001B[0;31mNameError\u001B[0m                                 Traceback (most recent call last)",
   "Cell \u001B[0;32mIn [2], line 1\u001B[0m\n\u001B[0;32m----> 1\u001B[0m \u001B[43mjob_input\u001B[49m\u001B[38;5;241m.\u001B[39mexecute_query(\u001B[38;5;124m\"\u001B[39m\u001B[3
   "\u001B[0;31mNameError\u001B[0m: name 'job_input' is not defined"
  ]
 }
],
"source": [
 "job_input.execute_query(\"DROP TABLE IF EXISTS rest_target_table;\")"
]
```

VERSATILE
DATA KIT

# Version Control Challenges

## Without VDK

```
350        {
351          "cell_type": "code",
  -          "execution_count": null,
352  +        "execution_count": 4,
353          "id": "cc05260f-2457-4174-9788-f185b24dd821",
354          "metadata": {
355            "tags": [
356              "vdk"
357            ]
358          },
  -          "outputs": [],
359  +        "outputs": [
360  +          {
361  +            "ename": "NameError",
362  +            "evalue": "name 'job_input' is not defined",
363  +            "output_type": "error",
364  +            "traceback": [
365  +              "\u001b[0;31m---------------------------------------------------------------------------\u001b[0m",
366  +              "\u001b[0;31mNameError\u001b[0m                                 Traceback (most recent call last)",
367  +              "Cell \u001b[0;32mIn [4], line 1\u001b[0m\n\u001b[0;32m----> 1\u001b[0m run(\u001b[43mjob_input\u001b[49m)\n",
368  +              "\u001b[0;31mNameError\u001b[0m: name 'job_input' is not defined"
369  +            ]
370  +          }
371  +        ],
372          "source": [
373            "run(job_input)"
374          ]
```

## With VDK

```
359        "outputs": [],
360        "source": [
361          "run(job_input)"
```

VERSATILE
DATA KIT

# From Data Exploration to Production



**PRODUCTION**

✓ Reproducibility: Non-Linear Execution and Hidden State Risks

✓ Code Organization: Irrelevant or debugging code

✓ Execution model: interactive kernel vs automated flow

✓ Automated Testing and CICD

✓ Version Control

**https://bit.ly/vdk-product-notebooks**

# Try it yourself

VERSATILE
DATA KIT

Please take the survey.

Thank You

https://www.linkedin.com/in/antoni-ivanov

VERSATILE
DATA KIT