# ⭐ Rules Based Enginee-2

Here is the updated code, including a utility for JSON validation and logging to avoid runtime failures caused by incorrect JSON values.

---

`feature_derivation_engine.py` **(Updated with Logging and Validation)**

```python
import logging
from pyspark.sql import functions as F

class FeatureDerivationEngine:
    def __init__(self, config):
        self.config = config
        self.logger = logging.getLogger("FeatureDerivationEngine")

    def validate_config(self):
        """
        Validates the configuration JSON to ensure it adheres to the expected schema.
        """
        required_keys = ["features"]
        for key in required_keys:
            if key not in self.config:
                raise ValueError(f"Missing required key '{key}' in configuration.")

        for feature in self.config["features"]:
            if "output_column" not in feature or "type" not in feature or "rule" not in feature:
                raise ValueError(f"Each feature must have 'output_column', 'type', and 'rule'.")
            if feature["type"] == "complex" and "conditions" not in feature["rule"]:
                raise ValueError(f"Complex type feature '{feature['output_column']}' must include 'conditions'.")
            if feature["type"] in ["transformation", "arithmetic"] and "expression" not in feature["rule"]:
                raise ValueError(f"Transformation type feature '{feature['output_column']}' must include 'expression'.")

        self.logger.info("Configuration validation successful.")

    def apply_complex_condition(self, df, conditions, default_value):
        exprs = None
        for condition in conditions:
            parent_condition = F.expr(condition["condition"])
            if "sub_conditions" in condition:
                sub_exprs = F.when(F.expr(condition["sub_conditions"][0]["condition"]),
    condition["sub_conditions"][0]["value"])
                for sub in condition["sub_conditions"][1:]:
                    sub_exprs = sub_exprs.when(F.expr(sub["condition"]), sub["value"])
                sub_exprs = sub_exprs.otherwise(default_value)
                parent_expr = F.when(parent_condition, sub_exprs)
            else:
                parent_expr = F.when(parent_condition, condition["value"])
            exprs = parent_expr if exprs is None else exprs.otherwise(parent_expr)
        exprs = exprs.otherwise(default_value)
        return exprs

    def apply_transformation(self, df, rule, output_col):
        expression = rule["expression"]
```

```
46            return df.withColumn(output_col, F.expr(expression))
47
48        def apply_rules(self, df):
49            self.logger.info("Starting feature derivation process.")
50            for feature in self.config["features"]:
51                output_col = feature["output_column"]
52                rule_type = feature["type"]
53                rule = feature["rule"]
54
55                self.logger.debug(f"Processing feature '{output_col}' of type '{rule_type}'.")
56
57                if rule_type == "complex":
58                    conditions = rule["conditions"]
59                    default_value = rule.get("default", "unknown")
60                    complex_expr = self.apply_complex_condition(df, conditions, default_value)
61                    df = df.withColumn(output_col, complex_expr)
62
63                elif rule_type in ["transformation", "arithmetic", "string"]:
64                    df = self.apply_transformation(df, rule, output_col)
65
66            self.logger.info("Feature derivation process completed.")
67            return df
68
```

### main.py (Updated with JSON Validation and Logging)

```
1   import json
2   import sys
3   import logging
4   from pyspark.sql import SparkSession
5   from src.feature_derivation_engine import FeatureDerivationEngine
6
7   def configure_logging():
8       logging.basicConfig(
9           level=logging.INFO,
10          format="%(asctime)s - %(name)s - %(levelname)s - %(message)s",
11          handlers=[
12              logging.StreamHandler(sys.stdout),
13          ],
14      )
15
16  def main(config_path):
17      # Configure logging
18      configure_logging()
19      logger = logging.getLogger("Main")
20
21      # Initialize Spark
22      spark = SparkSession.builder.appName("FeatureDerivation").getOrCreate()
23      logger.info("Spark session initialized.")
24
25      try:
26          # Load Configuration
27          logger.info("Loading configuration.")
28          with open(config_path, "r") as f:
29              config = json.load(f)
30
```

```
31            # Initialize Engine and Validate Config
32            engine = FeatureDerivationEngine(config)
33            logger.info("Validating configuration.")
34            engine.validate_config()
35
36            # Sample Data
37            data = [
38                (1, 15, 5), (2, 15, 15), (3, 15, 25), (4, 30, 50), (5, 70, 50)
39            ]
40            columns = ["id", "age", "score"]
41            df = spark.createDataFrame(data, columns)
42
43            # Apply Feature Derivation
44            logger.info("Applying feature derivation.")
45            result_df = engine.apply_rules(df)
46
47            # Show Results
48            logger.info("Feature derivation completed. Showing results:")
49            result_df.show()
50
51        except Exception as e:
52            logger.error(f"An error occurred: {e}")
53            sys.exit(1)
54
55  if __name__ == "__main__":
56        if len(sys.argv) != 2:
57            print("Usage: spark-submit main.py <config_file>")
58            sys.exit(-1)
59
60        config_path = sys.argv[1]
61        main(config_path)
62
```

### test_feature_derivation_engine.py (Updated with Validation Tests)

```
1   import unittest
2   from pyspark.sql import SparkSession
3   from src.feature_derivation_engine import FeatureDerivationEngine
4
5   class FeatureDerivationTest(unittest.TestCase):
6       @classmethod
7       def setUpClass(cls):
8           cls.spark = SparkSession.builder.appName("UnitTest").getOrCreate()
9
10      @classmethod
11      def tearDownClass(cls):
12          cls.spark.stop()
13
14      def test_valid_config(self):
15          config = {
16              "features": [
17                  {"output_column": "category", "type": "complex", "rule": {"conditions": [{"condition": "age <
    18", "value": "minor"}]}}
18              ]
19          }
20          engine = FeatureDerivationEngine(config)
```

```
21          try:
22              engine.validate_config()
23          except Exception as e:
24              self.fail(f"Validation failed unexpectedly: {e}")
25
26      def test_invalid_config_missing_keys(self):
27          config = {"features": [{"type": "complex", "rule": {"conditions": [{"condition": "age < 18", "value":
    "minor"}]}}]}
28          engine = FeatureDerivationEngine(config)
29          with self.assertRaises(ValueError):
30              engine.validate_config()
31
32      def test_apply_rules(self):
33          config = {
34              "features": [
35                  {"output_column": "category", "type": "complex", "rule": {"conditions": [{"condition": "age <
    18", "value": "minor"}]}}
36              ]
37          }
38          data = [(1, 15), (2, 25)]
39          columns = ["id", "age"]
40          df = self.spark.createDataFrame(data, columns)
41
42          engine = FeatureDerivationEngine(config)
43          engine.validate_config()
44          result_df = engine.apply_rules(df)
45
46          expected_data = [(1, 15, "minor"), (2, 25, "unknown")]
47          expected_df = self.spark.createDataFrame(expected_data, columns + ["category"])
48          self.assertEqual(sorted(result_df.collect()), sorted(expected_df.collect()))
49
50  if __name__ == "__main__":
51      unittest.main()
52
```

### Key Enhancements

1. **Validation Utility**:
   - Ensures the JSON configuration adheres to the expected schema before runtime.
   - Catches configuration errors early.
2. **Logging**:
   - Adds detailed logging to track the execution flow and debug issues.
   - Logs initialization, validation, and rule application steps.

These enhancements make the application robust and easy to debug during failures or misconfigurations.

Here is the full Python project structure for the updated PySpark code, including instructions on how to compile and package it into a deployable package.

## Project Directory Structure

```
1  feature-derivation-project/
```

```
 2  ├── config/
 3  │   └── sample_config.json        # Sample JSON configuration file
 4  ├── src/
 5  │   ├── __init__.py
 6  │   └── feature_derivation_engine.py # Feature Derivation Engine implementation
 7  ├── tests/
 8  │   ├── __init__.py
 9  │   └── test_feature_derivation_engine.py # Unit tests
10  ├── main.py                        # Entry point for the PySpark job
11  ├── requirements.txt               # Python dependencies
12  ├── setup.py                       # Packaging script
13  └── README.md                      # Documentation
14
```

## File Contents

**1.** `config/sample_config.json`

```json
 1  {
 2    "features": [
 3      {
 4        "output_column": "category",
 5        "type": "complex",
 6        "rule": {
 7          "conditions": [
 8            { "condition": "age < 18", "value": "child" },
 9            { "condition": "age >= 18 AND age < 60", "value": "adult" },
10            { "condition": "age >= 60", "value": "senior" }
11          ],
12          "default": "unknown"
13        }
14      }
15    ]
16  }
17
```

**2.** `src/__init__.py`

```
 1  # Empty init file to mark this as a package
 2
```

**3.** `src/feature_derivation_engine.py`

(Refer to the updated implementation provided earlier.)

**4.** `tests/__init__.py`

```
 1  # Empty init file to mark this as a package
 2
```

**5.** `tests/test_feature_derivation_engine.py`

(Refer to the updated test implementation provided earlier.)

**6.** `main.py`

(Refer to the updated implementation provided earlier.)

**7.** `requirements.txt`

```
1  pyspark==3.4.0
2
```

**8.** `setup.py`

```python
1   from setuptools import setup, find_packages
2
3   setup(
4       name="feature_derivation_project",
5       version="1.0.0",
6       description="Feature Derivation Engine for PySpark",
7       author="Your Name",
8       packages=find_packages(),
9       install_requires=[
10          "pyspark==3.4.0",
11      ],
12      entry_points={
13          "console_scripts": [
14              "feature-derivation=main:main"
15          ]
16      },
17  )
18
```

**9.** `README.md`

```
1   # Feature Derivation Project
2
3   This project implements a feature derivation engine in PySpark, supporting configurable and rule-based
    transformations.
4
5   ## How to Use
6
7   ### Configuration
8   Provide a JSON configuration file specifying the features and rules. Example:
9
10  ```json
11  {
12    "features": [
13      {
14        "output_column": "category",
15        "type": "complex",
16        "rule": {
17          "conditions": [
18            { "condition": "age < 18", "value": "child" },
19            { "condition": "age >= 18 AND age < 60", "value": "adult" },
20            { "condition": "age >= 60", "value": "senior" }
21          ],
22          "default": "unknown"
23        }
24      }
25    ]
26  }
27
```

**Run the Application**

Run the project using `spark-submit`:

```
1  spark-submit main.py config/sample_config.json
2
```

**Testing**

Run unit tests with `unittest`:

```
1  python -m unittest discover -s tests
2
```

**Packaging**

To build a deployable package:

```
1  python setup.py sdist
2
```

Deploy the generated `.tar.gz` file as needed.

```
1
2  ---
3
4  ## **How to Compile and Package**
5
6  ### 1. **Install Dependencies**
7  Set up a virtual environment and install dependencies:
8
9  ```bash
10 python -m venv venv
11 source venv/bin/activate
12 pip install -r requirements.txt
13
```

**2. Run Tests**

Run the unit tests to ensure everything works:

```
1  python -m unittest discover -s tests
2
```

**3. Package the Project**

Package the project into a `.tar.gz` file:

```
1  python setup.py sdist
2
```

This will generate a `dist/` directory containing the package file, e.g., `feature_derivation_project-1.0.0.tar.gz`.

**4. Deploy and Run**

Use `spark-submit` to run the packaged project:

```
1  spark-submit --py-files dist/feature_derivation_project-1.0.0.tar.gz main.py config/sample_config.json
2
```

**Logging Notes**

The logging information is displayed on the console and can be configured to write to a file if needed. Update the `logging.basicConfig()` in `main.py` for file logging.

This setup ensures your project is modular, testable, and deployable!