# Contents

# Validation Capabilities

1. **Basic Value Checks:**
   - Check if a column contains only unique values.
   - Verify if a column has no missing values.
   - Validate if a column contains only certain expected values.
2. **Value Range Checks:**
   - Ensure that numeric columns fall within specified ranges.
   - Validate dates fall within a certain timeframe.
   - Check if values in a column are within certain bounds.
3. **Pattern and Format Checks:**
   - Validate email addresses using regular expressions.
   - Check if strings follow specific formatting rules.
   - Ensure consistent casing or capitalization.
4. **Distribution Checks:**
   - Verify that numeric columns follow a specific distribution (e.g., normal distribution).
   - Check if categorical values are distributed as expected.
5. **Relationship Checks:**
   - Validate relationships between columns (e.g., total amount matches sum of individual items).
   - Ensure referential integrity across related tables.
6. **Completeness Checks:**
   - Validate that all required columns are present in the dataset.
   - Verify that specific combinations of columns are present.
7. **Uniqueness Checks:**
   - Check for duplicated rows or unique key constraints.
   - Ensure that combinations of columns are unique.
8. **Row Count and Summary Checks:**
   - Validate the total number of rows in a dataset.
   - Perform summary statistics checks (e.g., mean, median, standard deviation).
9. **Schema Checks:**
   - Validate that the data adheres to an expected schema (e.g., column names, data types).
10. **Custom Validations:**
    - Implement custom validation functions using Python code.
    - Combine multiple checks to perform complex validations.
11. **Time-based Validations:**
    - Ensure time-based data is in chronological order.
    - Validate that time intervals are consistent.
12. **Aggregation Validations:**

- Validate aggregate values against expected values (e.g., sum of sales matches a predefined value).
13. **File and Data Source Checks:**
    - Validate data read from files, databases, or external APIs.
    - Check if the data source schema matches expectations.
14. **Data Drift Monitoring:**
    - Compare data distributions over time to identify changes.
15. **Missing Data Checks:**
    - Validate the percentage of missing data in columns.
16. **Advanced Machine Learning Checks:**
    - Check if model predictions match expected outcomes.

These are just a few examples of the wide range of data validations you can perform using Great Expectations. The library is highly customizable and can accommodate a variety of use cases for ensuring data quality and integrity throughout your data pipeline.

# Pros

- **Comprehensive Validation:** Covers a wide range of data quality checks and complex validations.
- **Automated Testing:** Automates validation processes, reducing errors and ensuring consistent data quality.
- **Documentation Generation:** Automatically creates detailed documentation for expectations and validation results.
- **Collaboration:** Supports collaboration between data teams and domain experts.
- **Custom Expectations:** Enables creation of tailored expectations for specific use cases.
- **Continuous Monitoring:** Facilitates ongoing data quality monitoring.
- **Integration:** Can be integrated into various data processing environments.

# Cons

- **Setup Complexity:** Initial setup and expectation definition require learning and effort.
- **Maintenance:** Expectations need updates as data changes, leading to ongoing maintenance.
- **Resource Intensive:** Complex validations on large datasets might require significant resources.
- **Limited Transformation:** Not designed for extensive data transformation.
- **Configuration Dependency:** Accuracy depends on correct expectation configuration.
- **Learning Curve for Customization:** Customizing expectations requires additional learning.
- **Complexity for Small Projects:** Overhead might outweigh benefits for smaller projects.

# Conclusion

**Best For:**

Great Expectations excels in **complex data pipeline validation**, especially when you need to validate data quality and consistency at **multiple stages of a data pipeline**. It's a robust choice for data validation and documentation in data engineering workflows.

**Strengths:**

Great Expectations offers **comprehensive validation features**, automated documentation generation, and integration with various data processing tools. It's particularly powerful for **data monitoring** and validation in data pipelines.

**Considerations:**

Great Expectations **might be overkill for simpler validation** tasks or when you're primarily working with pandas DataFrames in a data analysis context.

## PANDERA

# Validation Capabilities

1. **Data Types:**
   - Check if columns have the expected data types (e.g., integers, floats, strings).
2. **Value Checks:**
   - Ensure values meet certain conditions (e.g., positive numbers, non-negative integers).
3. **Regular Expressions:**
   - Validate string columns against specific regular expressions.
4. **String Checks:**
   - Verify string values meet certain criteria (e.g., length constraints, specific substrings).
5. **Uniqueness Checks:**
   - Confirm that values in a column are unique (or not unique) within the DataFrame.
6. **Missing Data:**
   - Validate the presence or absence of missing values (NaN or None).
7. **Range Checks:**
   - Ensure that numeric values fall within specified ranges (e.g., between 0 and 100).
8. **Set Membership:**
   - Check if values belong to specific sets (e.g., categories, predefined sets of values).
9. **Date and Time Checks:**
   - Validate date and time columns (e.g., date format, future dates, past dates).
10. **Aggregate Functions:**
    - Apply aggregate functions (e.g., sum, mean, median) and validate their results.
11. **Custom Validation Functions:**
    - Define custom validation functions to check specific conditions or complex criteria.
12. **Column Existence:**
    - Ensure that certain columns exist in the DataFrame.

13. **Data Frame Shape:**
    - Validate the number of rows and columns in the DataFrame.
14. **Schema Validation:**
    - Check if the DataFrame adheres to an expected schema, including column names and data types.
15. **Column Dependence:**
    - Validate dependencies between columns (e.g., if one column has a certain value, another column should have a specific value).
16. **Multi-column Validation:**
    - Perform validations that involve multiple columns (e.g., verifying that a start date is before an end date).
17. **Hierarchical Data:**
    - Validate hierarchical or nested data structures within a DataFrame.
18. **Data Drift Detection:**
    - Detect changes in data distributions or statistics between different time periods.
19. **Combining Validations:**
    - Combine multiple validation rules for more complex checks.
20. **Data Transformation:**
    - Apply data transformations to the DataFrame and validate the transformed data.
21. **Integration with External Services:**
    - Integrate Pandera with external data validation services or APIs.

Pandera's flexibility allows to define comprehensive data validation rules to ensure data quality and integrity in pandas DataFrames.

# Pros

- **Integration with pandas:** Pandera is designed to work seamlessly with pandas DataFrames, making it well-suited for data validation in data analysis and manipulation workflows.
- **Declarative Syntax:** It offers a declarative and expressive syntax for defining validation rules, which enhances code readability and maintainability.
- **Comprehensive Validation:** Pandera supports a wide range of validation checks, from basic data type validation to complex custom checks and aggregations.
- **Schema Validation:** You can define and validate the entire schema of your Dataframe, including column names, data types, and validation rules, providing comprehensive data validation.
- **Customization:** You can easily define custom validation functions to accommodate specific data validation requirements.
- **Data Transformation:** Pandera allows you to define and validate data transformations, ensuring the quality of data processing steps.
- **Extensive Documentation:** The library has extensive documentation and examples, making it easier to get started and use effectively.

# Cons

- **Learning Curve:** While the library's core functionality is straightforward, mastering advanced features and custom checks may require some learning.
- **Setup Overhead:** There is an initial setup overhead in defining validation schemas, which can be cumbersome for very large or complex DataFrames.

- **Performance Considerations:** Applying complex validation rules to large DataFrames can introduce some performance overhead.
- **Python-Specific:** Pandera is specific to Python and pandas, limiting its use in multi-language or non-Python environments.
- **Maintenance:** Data validation schemas may need to be updated as data structures or validation requirements change.
- **Data Cleaning:** While Pandera validates data, it doesn't provide built-in data cleaning or transformation capabilities; this must be handled separately.
- **Limited to pandas DataFrames:** Pandera primarily targets pandas DataFrames, so it may not be the best choice for non-tabular data structures.

# Conclusion

**Best for:**

Pandera is best suited for data validation within **data analysis and pandas-based data manipulation** workflows. It's a good choice when you want to validate and ensure data quality within a pandas DataFrame.

**Strengths:**

Pandera provides an **expressive syntax for defining validation rules** and integrates seamlessly with pandas. It's well-suited for **declarative schema validatio**n and works well in data preprocessing pipelines.

**Considerations:**

While Pandera can be used for more complex validation scenarios, it's primarily **focused on tabular data,** and the learning curve might be steeper for those less familiar with pandas.

## PYDANTIC

# Validation Capabilities

Pydantic provides a wide range of validation capabilities beyond the basic data type checks. Here are some common and more advanced data validations you can perform using Pydantic:

1. **Data Types and Type Coercion:**
   - Check data types, such as `str`, `int`, `float`, `bool`, etc.
   - Use type coercion to automatically convert compatible data types.
2. **Field Constraints:**
   - Define constraints like minimum and maximum values for numeric types.
   - Define constraints like minimum and maximum lengths for strings.
3. **Regular Expressions:**
   - Use regular expressions to validate string formats, like email addresses, URLs, or custom patterns.
4. **Enums:**
   - Use Python enums to restrict a field's value to a predefined set of choices.
5. **Optional Fields:**
   - Use `Optional` to mark fields as optional, allowing them to be omitted in the input data.

6. **Default Values:**
   - Provide default values for fields that can be omitted in the input.
7. **Custom Validators:**
   - Define custom validation functions for fields. These functions should raise `ValueError` or return `None` for valid values and raise `ValueError` for invalid values.
8. **Conditional Validation:**
   - Conditionally validate fields based on the values of other fields. Use the `@validator` decorator for this purpose.
9. **Complex Data Structures:**
   - Validate more complex data structures like nested dictionaries and lists.
10. **List Length Validation:**
    - Check the length of lists and enforce minimum and maximum lengths.
11. **Dependency Validation:**
    - Perform validation based on the presence or absence of other fields.
12. **Validating JSON Data:**
    - Use Pydantic models to validate JSON data, including nested JSON objects.
13. **Validating Configuration Settings:**
    - Use Pydantic to validate and parse configuration settings from files or environment variables.
14. **Custom Error Messages:**
    - Customize error messages for validation failures by defining a `@root_validator` with `pre=True` and returning a `ValueError` with a custom message.
15. **Multi-Field Validation:**
    - Validate multiple fields together, for example, checking that a start date is before an end date.
16. **External Data Validation:**
    - Integrate with external validation services or databases to validate data against external sources.
17. **Validation Groups:**
    - Define validation groups to perform specific validations based on use cases or contexts.

These are just some examples of the validation capabilities Pydantic offers. Pydantic flexibility and integration with Python type hints make it a versatile tool for ensuring data quality and consistency in our applications. We can combine multiple validation techniques to handle complex validation scenarios effectively.

Here are some brief pros and cons of Pydantic library

# Pros

- **Type Hint Integration:** Uses Python type hints for clear and readable data structure definitions.
- **Data Validation:** Automatically validates input data based on defined types and constraints.
- **Type Coercion:** Can automatically convert compatible data types for you.
- **Default Values:** Supports setting default values for optional fields.
- **Custom Validation:** Allows custom validation logic beyond basic type checks.
- **Complex Structures:** Handles nested and complex data structures effectively.
- **Automated Docs:** Generates JSON Schema for API documentation.
- **Dependency Validation:** Validates data based on other field values.

# Cons

- **Learning Curve:** Advanced features might require time to master.

- **Performance Overhead:** Introduces a slight performance overhead for validation and type coercion.
- **Limited Transformation:** Primarily focused on validation, not complex transformations.
- **Complex Scenarios:** Intricate validation might need custom logic.
- **Code Clutter:** Annotations and validation can make code more complex.
- **Framework Integration:** Requires understanding specific integration with certain frameworks.
- **Custom Errors:** Customizing error messages can be challenging.
- **Python Only:** Limited to Python and not suitable for non-Python systems.

# Conclusion

**Best for:**

Pydantic is ideal for **data validation in Python applications**, including APIs, configuration settings, and data input validation. It's a versatile choice for enforcing **data quality** and ensuring **data consistency** in Python programs.

**Strengths:**

Pydantic leverages **Python's type hinting** system and integrates smoothly with Python code. It's particularly well-suited for validating and **parsing data before it's processed.**

**Considerations:**

While Pydantic can be used for data validation in various contexts, it may **not be as feature-rich as Great Expectations** for complex data pipeline validation scenarios.

## SUMMARY

# Pydantic Vs. Pandera Vs. Great Expectations

| Feature | Pandera | Pydantic | Great Expectations |
|---|---|---|---|
| **Integration** | Seamlessly integrates with pandas DataFrames | Integrates with Python applications | Designed for complex data pipelines |
| **Schema Validation** | Supports schema validation | Schema validation for Python types | Offers automated data testing and docs |
| **Custom Validation** | Allows custom validation functions | Supports custom validation functions | Supports custom validation functions |

| Feature | Pandera | Pydantic | Great Expectations |
|---|---|---|---|
| **Data Transformation** | Supports data transformation validation | Handles data parsing in Python apps | Primarily focuses on data engineering |
| **Complex Validation** | Suitable for data analysis workflows | Versatile for Python applications | Designed for multi-stage data pipelines |
| **Automated Testing** | Limited | Limited | Offers automated data testing |
| **Data Monitoring** | Limited | Limited | Provides continuous data quality checks |
| **Learning Curve** | Shallow, provides a familiar pandas like API, and Moderate for complex cases | Shallow/Moderate - Provides a straightforward Pythonic API. | Initial learning curve for setup is relatively steep |
| **Python Specific** | Yes (pandas integration) | Yes | Python-centric, may not suit non-Python |
| **Project Size Considerations** | Good for data analysis projects | Versatile for Python applications | Best suited for complex data pipelines |
| **Use with PySpark** | Can be used but requires DataFrame conversion to pandas for validation. Potential performance overhead with large DataFrames. | Use within Python code that interacts with PySpark for data validation. Validates data before or after PySpark processing. | Supports PySpark integration for end-to-end data validation within PySpark pipelines. Most suitable for comprehensive PySpark data validation |
| **Complex Hierarchy** | Limited support; primarily designed for tabular data. | Limited support; designed for validation in Python applications | Robust support for complex hierarchical data |

| Feature | Pandera | Pydantic | Great Expectations |
|---|---|---|---|
| **Structure Support** | | but not specialized for complex hierarchies. | structures, designed for data engineering and data pipeline scenarios. |
| **Documentation** | Excellent | Excellent | Good |
| **Writing Validation Tests Quickly** | Very good fit. Simple and familiar Python API makes this easy. Python code becomes the definition of the tests. | Not a great fit. Python API is less concise, and the Python code ultimately gets converted to JSON, which becomes the definition of the expectations. | Very good fit. Concise Pythonic syntax for defining and validating data models. |
| **Ability to Profile Data** | Yes. | Yes, and generates a more comprehensive set of tests. | Not explicitly designed for data profiling. |
| **Integration with Other Tools** | Features and documentation are relatively light on this front, suggesting the library is not well-suited to this requirement. | It's clear that integration has been a major consideration in the design of this library. | Can integrate with other Python libraries and services. |
| **Building Comprehensive Validation System in a Complex Domain** | Features and documentation are relatively light on this front, indicating that this library isn't well-suited to this requirement, or at least that fulfilling it would require considerable work. | It's clear that this library has been designed for the purpose of building comprehensive validation systems that can handle the complexities of the data space. | Suitable for data validation in Python applications; might require additional work for complex domain-specific validation systems. |