

# cuallee: A python package for data quality checks across multiple DataFrame APIs

Herminio Vazquez <sup>1\*</sup> and Virginie Grosboillot <sup>2\*</sup>

<sup>1</sup> Independent Researcher, Mexico <sup>2</sup> Swiss Federal Institute of Technology (ETH) \* These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

In today's world, where vast amounts of data are generated and collected daily, and where data heavily influence business, political, and societal decisions, it is crucial to evaluate the quality of the data used for analysis, decision-making, and reporting. This involves understanding how reliable and trustworthy the data are. To address this need, we have created cuallee, a Python package for assessing data quality. cuallee is designed to be dataframe-agnostic, offering an intuitive and user-friendly API for describing checks across the most popular dataframe implementations such as PySpark, Pandas, Snowpark, Polars, DuckDB and BigQuery. Currently, cuallee offers over 50 checks to help users evaluate the quality of their data.

## Statement of need

For data engineers and data scientists, maintaining a consistent workflow involves operating in hybrid environments, where they develop locally before transitioning data pipelines and analyses to cloud-based environments. Whilst working in local environments typically allows them to fit data sets in memory, moving workloads to cloud environments involve operating with full scale data that requires a different computing framework ([Schelter et al., 2018](#)), i.e. distributed computing, parallelization, and horizontal scaling. cuallee accomodates the testing activities required by this shift in computing frameworks, in both local and remote environments, without the need to rewrite test scenarios or employ different testing approaches for assessing various quality dimensions of the data ([Fadlallah et al., 2023b](#)).

An additional argument is related to the rapid evolution of the data ecosystem ([Fadlallah et al., 2023a](#)). Organizations and data teams are constantly seeking ways to improve, whether through cost-effective solutions or by integrating new capabilities into their data operations. However, this pursuit presents new challenges when migrating workloads from one technology to another. As information technology and data strategies become more resilient against vendor lock-ins, they turn to technologies that enable seamless operation across platforms, avoiding the chaos of fully re-implementing data products. In essence, with cuallee no data testing strategy needs to be rewritten or reformulated due to platform changes.

One last argument in favor of using a quality tool such as cuallee is the need to integrate quality procedures into the early stages of data product development. Whether in industry or academia, there is often a tendency to prioritize functional aspects over quality, leading to less time being dedicated to quality activities. By providing a clear, easy-to-use, and adaptable programming interface for data quality, teams can incorporate quality into their development process, promoting a proactive approach of building quality in rather than relying solely on testing to ensure quality.

40

## Data Quality Frameworks

41 Data platforms have diversified from file systems and relational databases, to full ecosystems  
42 including the concept of data lakes (Dumke et al., 2020). Modern platforms host a variety of  
43 data formats, beyond traditional tabular data, including semi-structured like JSON (Pezoa et  
44 al., 2016) or unstructured like audio or images.

45 Operating with modern data platforms, requires a sophisticated data processing framework  
46 capable to handle multiple formats, and scalability. Apache Spark (Armbrust et al., 2015)  
47 has revolutionized the data flow paradigm by bringing computation to the data, reversing the  
48 omnipresent data to the computation, it has commoditized large scale data processing and it  
49 has grown in adoption.

50 Apache Spark’s growth can be attributed to its ease of use, versatility, and performance. It  
51 supports multiple programming languages, including Python, Scala, Java, and R, making  
52 it accessible to a wide range of developers. Moreover, Spark’s ability to handle various  
53 data processing tasks —batch processing, real-time streaming, machine learning, and graph  
54 processing—within a unified framework has been a key factor in its widespread adoption  
55 (O’Reilly Media, 2023).

56 cuallee is powered by native data engines like Apache Spark. Compared to other data quality  
57 frameworks it brings substantial advantages in reduced complexity, less computation resources  
58 and the fastest time per validation.

59 The following table (Table 1) provides a summary of the performance comparisson:

Table 1: Performance comparisson on popular data quality frameworks

Framework	Checks Definition	Time
great_expectations	python	66s
soda	yaml	43s
pydeequ	python	11s
cuallee	python	7s

60

## Methods

61 cuallee employs a heuristic-based approach to define quality rules for each dataset. This  
62 prevents the inadvertent duplication of quality predicates, thus reducing the likelihood of  
63 human error in defining rules with identical predicates. Several studies have been conducted  
64 on the efficiency of these rules, including auto-validation and auto-definition using profilers  
65 (Tu et al., 2023).

66

## Checks

67 In cuallee, checks serve as the fundamental concept. These checks (Table 2) are implemented  
68 by rules, which specify quality predicates. These predicates, when aggregated, form the criteria  
69 used to evaluate the quality of a dataset. Efforts to establish a universal quality metric  
70 (Pleimling et al., 2022) typically involve using statistics and combining dimensions to derive a  
71 single reference value that encapsulates overall quality attributes.

**Table 2:** List and description of the currently available checks

Check	Description	Data Type
is_complete	Zero nulls	<i>agnostic</i>
is_unique	Zero duplicates	<i>agnostic</i>
is_primary_key	Zero duplicates	<i>agnostic</i>
are_complete	Zero nulls on group of columns	<i>agnostic</i>
are_unique	Composite primary key check	<i>agnostic</i>
is_composite_key	Zero duplicates on multiple columns	<i>agnostic</i>
is_greater_than	$col > x$	<i>numeric</i>
is_positive	$col > 0$	<i>numeric</i>
is_negative	$col < 0$	<i>numeric</i>
is_greater_or_equal_than	$col \geq x$	<i>numeric</i>
is_less_than	$col < x$	<i>numeric</i>
is_less_or_equal_than	$col \leq x$	<i>numeric</i>
is_equal_than	$col == x$	<i>numeric</i>
is_contained_in	$col \in [a, b, c, \dots]$	<i>agnostic</i>
is_in	Alias of is_contained_in	<i>agnostic</i>
not_contained_in	$col \notin [a, b, c, \dots]$	<i>agnostic</i>
not_in	Alias of not_contained_in	<i>agnostic</i>
is_between	$a \leq col \leq b$	<i>numeric, date</i>
has_pattern	Matching a pattern defined as a regex	<i>string</i>
is_legit	String not null & not empty $^{\wedge}\S^{\$}$	<i>string</i>
has_min	$\min(col) == x$	<i>numeric</i>
has_max	$\max(col) == x$	<i>numeric</i>
has_std	$\sigma(col) == x$	<i>numeric</i>
has_mean	$\mu(col) == x$	<i>numeric</i>
has_sum	$\Sigma(col) == x$	<i>numeric</i>
has_percentile	$\%(col) == x$	<i>numeric</i>
has_cardinality	$\text{count}(\text{distinct}(col)) == x$	<i>agnostic</i>
has_max_by	A utility predicate for $\max(col\_a) == x$ for $\max(col\_b)$	<i>agnostic</i>
has_min_by	A utility predicate for $\min(col\_a) == x$ for $\min(col\_b)$	<i>agnostic</i>
has_correlation	Finds correlation between 0..1 on $\text{corr}(col\_a, col\_b)$	<i>numeric</i>
has_entropy	Calculates the entropy of a column $\text{entropy}(col) == x$ for classification problems	<i>numeric</i>
is_inside_iqr	Verifies column values reside inside limits of interquartile range $Q1 \leq col \leq Q3$ used on anomalies.	<i>numeric</i>
is_in_millions	$col \geq 1e6$	<i>numeric</i>
is_in_billions	$col \geq 1e9$	<i>numeric</i>
is_t_minus_1	For date fields confirms 1 day ago $t-1$	<i>date</i>
is_t_minus_2	For date fields confirms 2 days ago $t-2$	<i>date</i>
is_t_minus_3	For date fields confirms 3 days ago $t-3$	<i>date</i>
is_t_minus_n	For date fields confirms n days ago $t-n$	<i>date</i>
is_today	For date fields confirms day is current date $t-0$	<i>date</i>
is_yesterday	For date fields confirms 1 day ago $t-1$	<i>date</i>
is_on_weekday	For date fields confirms day is between Mon-Fri	<i>date</i>

Check	Description	Data Type
is_on_weekend	For date fields confirms day is between Sat-Sun	<i>date</i>
is_on_monday	For date fields confirms day is Mon	<i>date</i>
is_on_tuesday	For date fields confirms day is Tue	<i>date</i>
is_on_wednesday	For date fields confirms day is Wed	<i>date</i>
is_on_thursday	For date fields confirms day is Thu	<i>date</i>
is_on_friday	For date fields confirms day is Fri	<i>date</i>
is_on_saturday	For date fields confirms day is Sat	<i>date</i>
is_on_sunday	For date fields confirms day is Sun	<i>date</i>
is_on_schedule	For date fields confirms time windows i.e. 9:00 - 17:00	<i>timestamp</i>
is_daily	Can verify daily continuity on date fields by default. [2,3,4,5,6] which represents Mon-Fri in PySpark. However new schedules can be used for custom date continuity	<i>date</i>
has_workflow	Adjacency matrix validation on 3-column graph, based on group, event, order columns.	<i>agnostic</i>
satisfies	An open SQL expression builder to construct custom checks	<i>agnostic</i>
validate	The ultimate transformation of a check with a dataframe input for validation	<i>agnostic</i>
iso.iso_4217	currency compliant ccy	<i>string</i>
iso.iso_3166	country compliant country	<i>string</i>
Control.completeness	Zero nulls all columns	<i>agnostic</i>
Control.percentage_fill	% rows not empty	<i>agnostic</i>
Control.percentage_empty	% rows empty	<i>agnostic</i>

## References

- Armbrust, M., Xin, R. S., Lian, C., Huai, Y., Liu, D., Bradley, J. K., Meng, X., Kaftan, T., Franklin, M. J., Ghodsi, A., & Zaharia, M. (2015). Spark SQL: Relational data processing in spark. *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, 1383–1394. <https://doi.org/10.1145/2723372.2742797>
- Dumke, A. R., Parchmann, A., Schmid, S., & Hauswirth, M. (2020). Toward data lakes as central building blocks for data management and analysis. *Frontiers in Big Data*, 3, 564115. <https://doi.org/10.3389/fdata.2020.564115>
- Fadlallah, H., Kilany, R., Dhayne, H., El Haddad, R., Haque, R., Taher, Y., & Jaber, A. (2023a). BIGQA: Declarative big data quality assessment. *Journal of Data and Information Quality*, 15. <https://doi.org/10.1145/3603706>
- Fadlallah, H., Kilany, R., Dhayne, H., El Haddad, R., Haque, R., Taher, Y., & Jaber, A. (2023b). Context-aware big data quality assessment: A scoping review. *Journal of Data and Information Quality*, 15. <https://doi.org/10.1145/3603707>
- O'Reilly Media. (2023). *Technology trends for 2023*. <https://www.oreilly.com/radar/technology-trends-for-2023/>
- Pezoa, F., Reutter, J., Suarez, F., Ugarte, M., & Vrgoč, D. (2016). *Foundations of JSON schema*. 263–273. <https://doi.org/10.1145/2872427.2883029>

- 90 Pleimling, X., Shah, V., & Lourentzou, I. (2022). [Data] quality lies in the eyes of the beholder.  
91 *Proceedings of the 15th International Conference on Pervasive Technologies Related to*  
92 *Assistive Environments*, 118–124. <https://doi.org/10.1145/3529190.3529222>
- 93 Schelter, S., Lange, D., Schmidt, P., Celikel, M., Biessmann, F., & Grafberger, A. (2018).  
94 Automating large-scale data quality verification. *Proc. VLDB Endow.*, 11(12), 1781–1794.  
95 <https://doi.org/10.14778/3229863.3229867>
- 96 Tu, D., He, Y., Cui, W., Ge, S., Zhang, H., Han, S., Zhang, D., & Chaudhuri, S. (2023).  
97 Auto-validate by-history: Auto-program data quality constraints to validate recurring data  
98 pipelines. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and*  
99 *Data Mining*, 4991–5003. <https://doi.org/10.1145/3580305.3599776>

DRAFT