

# PYRA: A High-level Linter for Data Science Software

Greta Dolcetti<sup>a</sup>, Vincenzo Arceri<sup>b</sup>, Antonella Mensi<sup>c</sup>, Enea Zaffanella<sup>b</sup>, Caterina Urban<sup>d</sup> and Agostino Cortesi<sup>a</sup>

<sup>a</sup>*Ca' Foscari University of Venice, Via Torino, 155, Venice, 30170, Italy*

<sup>b</sup>*University of Parma, Parco Area delle Scienze, 53/A, Parma, 43124, Italy*

<sup>c</sup>*University of Verona, Piazzale L. A. Scuro, 10, Verona, 37134, Italy*

<sup>d</sup>*Inria & École Normale Supérieure | Université PSL, Paris, France*

## ARTICLE INFO

### Keywords:

Static Analysis  
Jupyter Notebooks  
Data Science

## ABSTRACT

Due to its interdisciplinary nature, the development of data science software is particularly prone to a wide range of potential mistakes that can easily and silently compromise the final results. Several tools have been proposed that can help the data scientist in identifying the most common, low-level programming issues. However, these tools often fall short in detecting higher-level, domain-specific issues typical of data science pipelines, where subtle errors may not trigger exceptions but can still lead to incorrect or misleading outcomes, or unexpected behaviors.

In this paper, we present PYRA, a static analysis tool that aims at detecting code smells in data science workflows. PYRA builds upon the Abstract Interpretation framework to infer abstract datatypes, and exploits such information to flag 16 categories of potential code smells concerning misleading visualizations, challenges for reproducibility, as well as misleading, unreliable or unexpected results. Unlike traditional linters, which focus on syntactic or stylistic issues, PYRA reasons over a domain-specific type system to identify data science-specific problems – such as improper data preprocessing steps and procedures’ misapplications – that could silently propagate through a data-manipulation pipeline. Beyond static checking, we envision tools like PYRA becoming integral components of the development loop, with analysis reports guiding correction and helping assess the reliability of machine learning pipelines. We evaluate PYRA on a benchmark suite of real-world Jupyter notebooks, showing its effectiveness in detecting practical data science issues, thereby enhancing transparency, correctness, and reproducibility in data science software.

## 1. Introduction

Data science informally refers to an interdisciplinary field that integrates concepts from statistics, informatics, computing, communication, management, and sociology to analyze data and its environment (including domain-specific, organizational, and societal aspects). The ultimate aim of this discipline is to extract valuable insights from data that can be used for interpretative purposes or to assist in decision-making, following a data-to-knowledge-to-wisdom approach and methodology [3]. Given the widespread adoption of data science-based approaches across various fields – healthcare, retail, manufacturing, finance, etc. – several data science tools and libraries have become widely popular. These include, but are not limited to:

- scikit-learn [27], a Python library that allows the development of a complete machine learning pipeline;
- pandas [20], a Python library for data manipulation and analysis;

- seaborn [48] and ggplot2 [49], which are data visualization tools designed for Python and R, respectively;
- Jupyter Notebooks [16], a web application that, through the use of notebooks, allows to write and execute code, visualize data and add comments within one interface;
- BioConductor [11], an R ecosystem that encompasses a wide variety of bioinformatic tools.

This list of tools and libraries also shows that Python and R are the programming languages of choice for data scientists. Both languages are dynamically typed, meaning that they perform their type correctness checks at runtime and do not enforce native support for a more systematic, static control of the operations that are allowed on the values of variables; this means that a typing error in a seldomly executed computational path will only be discovered when running a test that actually triggers the execution of that specific computational path. In contrast, statically typed languages perform most (sometimes all) of the type checks before running the program, checking all its possible execution paths: hence, they can eagerly spot the most common programming errors even before running a single dynamic test.

It is worth stressing that the mere adoption a statically typed language would provide no guarantee on the code being completely correct: the type checking tool (typically run as a step in the compilation phase) will spot all proper typing errors, but logical errors would remain undetected; when present, logical errors can lead to unwanted or misleading results that the user may wrongly accept as correct.

\*Corresponding author

✉ greta.dolcetti@unive.it (G. Dolcetti); vincenzo.arceri@unipr.it (V. Arceri); antonella.mensi@univr.it (A. Mensi); enea.zaffanella@unipr.it (E. Zaffanella); caterina.urban@inria.fr (C. Urban); cortesi@unive.it (A. Cortesi)

ORCID(s): 0000-0002-2983-9251 (G. Dolcetti); 0000-0002-5150-0393 (V. Arceri); 0000-0001-9468-5298 (A. Mensi); 0000-0001-6388-2053 (E. Zaffanella); 0000-0002-8127-9642 (C. Urban); 0000-0002-0946-5440 (A. Cortesi)

47 Experience has shown that a significant percentage of these  
 48 logical errors can still be related to the “data type” of the  
 49 program variables, provided the default type system of the  
 50 considered programming language is replaced by a non-  
 51 standard, higher level type system, suitably extended so as  
 52 to detect and propagate the relevant information. For these  
 53 scenarios, several *ad hoc* type systems have been developed:  
 54 for instance, *session types* have been developed to help in  
 55 checking that a concurrent program fulfills the requirements  
 56 of a given communication protocol [10]; in safety critical  
 57 contexts, the MISRA-C coding standard [21] defines the  
 58 *essential type system* (among other things forbidding some  
 59 of the implicit type conversions that are legal for C code) and  
 60 requires that the program is well typed according to its rules.

61 The approaches above have in common the fact that  
 62 these non-standard type systems have a *prescriptive* nature: a  
 63 deviation from the typing rules is considered an error which  
 64 should be corrected. However, such a clear-cut distinction  
 65 between correct and wrong code cannot always be made.  
 66 In the cases where the tool identifies a *smell* in the code  
 67 the prescriptive approach is better replaced by a *descriptive*  
 68 approach, where the tool stops pretending to have a complete  
 69 knowledge and does its best to help the developer in under-  
 70 standing what is going on. For instance, almost all compilers  
 71 can issue a rich set of warnings: when clear and to the  
 72 point, this feedback is useful and greatly appreciated by the  
 73 programmer. This is also the reason for the development of  
 74 *linter* tools, i.e., lightweight tools that assist the programmer  
 75 in improving code quality by spotting questionable code.  
 76 Available linter tools differ in two main dimensions: the  
 77 considered programming language and the kind of issues  
 78 they focus on. The latter ranges from low level issues (e.g.,  
 79 respecting variable naming conventions or software metric  
 80 thresholds) to higher level issues, which often take into  
 81 account the intended semantics of a portion of code.

82 A proposal for the development of a linter tool for  
 83 data science code, focused on the Python language, was  
 84 put forward in [8]. The tool aims at detecting several data  
 85 science related code smells by gathering information about  
 86 the potential runtime values of variables into an *abstract type*  
 87 *system*. The latter comprises high-level data types tailored  
 88 specifically for data science code. Lastly, the tool verifies  
 89 that calls to data science library functions are consistent  
 90 with the determined abstract data types. As explained above,  
 91 the tool adopts a *descriptive* approach: its end goal is to  
 92 make the user reason about their code by reporting them  
 93 a list of putative inappropriate behaviors, without obliging  
 94 them to take a specific action; this fits rather well with  
 95 the fact that data science code is highly context-dependent.  
 96 The usefulness of this prototype is further enhanced by the  
 97 fact that many data scientists are not code specialists, e.g.,  
 98 software engineers or professional developers. Indeed, data  
 99 science is interdisciplinary, and the tools we have mentioned,  
 100 such as pandas, are highly user-friendly for anyone with a  
 101 basic understanding of programming.

102 In this paper we thoroughly extend [8, 7] and we present  
 103 PYRA, a working prototype of the linter tool that is easy

104 to use and integrates seamlessly with Python code, with-  
 105 out requiring additional annotations or modifications of the  
 106 code. The abstract datatype domain of PYRA comprises 56  
 107 datatypes – ranging from higher-level ones to others that are  
 108 data science-specific – designed to capture 16 categories of  
 109 the most common code smells, of various nature and gravity.

110 The implementation of PYRA is based on LYRA [44], a  
 111 static analyzer for Python that automatically detects input  
 112 data that remains unused by a Python program. It is a re-  
 113 search prototype and its support for Jupyter notebook is only  
 114 a proof of concept. It does not support any other detection of  
 115 domain-specific issues as PYRA. More concretely, [8] lays  
 116 the foundations for PYRA by motivating the need for a linter  
 117 for data science code: the notion of code smells specific to  
 118 data science is introduced using minimal examples, while  
 119 formally describing the adopted abstract domain and the  
 120 corresponding type rules. A refined version of the prototype  
 121 introduced in [8] is informally presented in [7], where its  
 122 functionalities and its utility are demonstrated adopting a  
 123 more practical point of view.

124 Building upon the previous work, in this paper we de-  
 125 scribe a further improved version of the tool, characterized  
 126 by additional checkers and a more robust implementation;  
 127 the contributions also include a more detailed description of  
 128 the tool’s behavior, with an explanation and classification of  
 129 the warnings produced, as well as an experimental evalua-  
 130 tion conducted on real notebooks, resulting in a significant  
 131 advancement compared to earlier efforts. We argue that  
 132 the Abstract Interpretation framework [5], due to its ability  
 133 to formalize approximation and support abstract domain  
 134 refinement, is particularly well-suited for the incremental  
 135 development of a descriptive (i.e., permissive) type system.

136 The rest of the paper is organized as follows. In Section 2  
 137 we briefly cover the related work, whereas in Section 3  
 138 we provide an overview on the code smells that we aim to  
 139 detect, categorize them and describe some of them in detail.  
 140 Section 4 thoroughly describes the proposed tool, PYRA,  
 141 covering its architecture, its abstract datatype domain, the  
 142 implemented checkers, and an example of its execution.  
 143 Lastly, Section 5 is dedicated to the experimental evaluation,  
 144 Section 6 discusses some limitations and important notes  
 145 and in Section 7 we draw some conclusions and discuss  
 146 potential ideas for future research.

## 2. Related Work

147 Abstract Interpretation [5] is a mathematical framework  
 148 that allows to formally derive approximations of the seman-  
 149 tics of programming languages. Its most common applica-  
 150 tion is the systematic development of sound static analyzers,  
 151 i.e., tools that are able to automatically infer some properties  
 152 of a program without executing it. In particular, [4] shows  
 153 how type systems and type inference algorithms can be cast  
 154 as instances of Abstract Interpretation. A gentle introduction  
 155 to the modeling of simple type information as Abstract Inter-  
 156 pretation is the *dimension calculus* of [6, Section 2.2]: here  
 157 it is shown how concrete unit of measures (e.g., meter, yard,  
 158

159 second, hour, kilogram, pound, ...) can be approximated  
 160 using abstract dimensions (e.g., length, time, mass, surface,  
 161 speed, ...) and then propagated via abstract rules such as

$$\begin{aligned} \text{length} + \text{length} &= \text{length}, \\ \text{length} \times \text{length} &= \text{surface}, \\ \text{length} / \text{length} &= \text{nodimension}, \\ \text{length} / \text{time} &= \text{speed}, \\ &\dots \end{aligned}$$

162 This simple idea can be easily generalized to more sophisticated  
 163 type systems, such as the one we propose in this paper.

164 Due to the importance and pervasiveness of data science,  
 165 the need to analyze Jupyter Notebooks has been highlighted [47], and many techniques to analyze data sciences  
 166 code have been proposed accordingly. For example, [24,  
 167 42, 43] propose a framework based on Abstract Interpretation [5] to infer necessary conditions on the structure  
 168 and values of the data read by a data-processing program  
 169 or to automatically detect unused input data [44]. Other  
 170 static analysis frameworks focus on detecting data leakage [9, 38, 39] or studying the impact of code changes across  
 171 code cells in notebooks. On the other end, open-source  
 172 tools like pandera [1] and pymlint [28] have been released  
 173 with the aim to perform data validation using schemas (i.e.  
 174 the specification of the expected structure, data types and  
 175 validation rules for the data), and reveal potential notebook  
 176 defects, recommending corrective actions that promote best  
 177 practices such as using version control and putting import  
 178 statements at the beginning of the notebook. Regarding static  
 179 type analysis and inference, many tools based on Abstract  
 180 Interpretation, such as [17, 22], or relying on Z3 [23] or  
 181 other SMT solvers, such as [13], have been proposed. How-  
 182 ever, these tools typically focus on inferring Python type  
 183 hints [30] and detecting potential errors. They usually target  
 184 the standard Python language and some standard libraries  
 185 (e.g., os, json), aiming to infer concrete type hints and errors.  
 186 In contrast, our goal is to infer and reason about more  
 187 abstract datatypes, potentially capturing a broader and less  
 188 conventional set of errors and code smells. Our work is  
 189 inspired by these projects but aims at finding more subtle  
 190 code smells and proposing an easily extensible framework  
 191 to help developers achieve correct results.

192 Even though not strictly related to the analysis of Jupyter  
 193 notebooks, research on the R programming language, another  
 194 one of the most popular languages for data and statistical  
 195 analysis, is also noteworthy. In [35], the authors  
 196 conducted a large-scale analysis of R programs, considering  
 197 both scripts submitted with academic publications and those  
 198 found in CRAN packages, investigating the most popular  
 199 features, constructs and operations of R. Based on this  
 200 study, [36] proposed flowR, a static dataflow analyzer and  
 201 program slicer for R programs, which also supports its  
 202 most challenging features, such as redefinition of primitive  
 203 constructs. Finally, in [12], the authors propose a large-scale  
 204 study on the usage of eval in R. They demonstrate that R  
 205 allows a higher degree of flexibility in using eval compared

206 to JavaScript, and they discuss the challenges associated  
 207 with analyzing or refactoring code that employs eval while  
 208 preserving its intended semantics.

209 To the best of our knowledge, there is not another frame-  
 210 work specifically designed to infer and reason about abstract  
 211 datatypes in Jupyter Notebooks and to capture a variety  
 212 of data science code smells by also using concrete dataset  
 213 information, as we do in PYRA. The most similar framework  
 214 is MLScents [32], even though it focuses on lower level anti-  
 215 patterns detection (e.g. missing docstring for function, magic  
 216 numbers, array creation efficiency, etc.) and it only uses a  
 217 fully static abstract syntax tree analysis. However, as shown  
 218 in Section 5, on the two issues that can be detected by both  
 219 tools, PYRA outperforms MLScents. Therefore, we claim that  
 220 PYRA is the first framework that combines Abstract Inter-  
 221 pretation with concrete dataset information to infer abstract  
 222 datatypes and detect a wide range of data science code smells  
 223 in Jupyter Notebooks.

### 3. Code Smells

224 In this section we provide an informal definition for what  
 225 we call a *data science code smell*, along with the issues  
 226 related to them and some minimal examples.

227 Generally speaking, a code smell is any characteristics  
 228 of (a portion of) the source code that hints at the existence  
 229 of a deeper problem, thereby hindering software mainte-  
 230 nance and evolution [26]. Even though code smells are  
 231 not necessarily bugs, they might cause issues and usually  
 232 denote a weakness in the code design. In the context of data  
 233 science code, we refine the definition above to mean any  
 234 code denoting an operation that, while being legal according  
 235 to the language of choice (i.e., it has a well defined behavior  
 236 and does not raise an exception), it may be a logical or  
 237 methodological mistake, potentially leading to computing  
 238 results that are incorrect in the considered context.

239 As mentioned in Section 1, PYRA focuses on code smells  
 240 that are specific to the data science pipeline when using the  
 241 Python language. The set of 16 categories of code smells  
 242 analyzed by PYRA was constructed by considering some  
 243 of the most common and well-known issues that can arise  
 244 in data science pipelines [50, 31, 18, 15], as well as some  
 245 other general issues that can lead to misleading results or  
 246 unexpected behaviors.

247 In this section, we provide descriptions and examples of  
 248 the most representative ones, while a brief overview of all  
 249 the included issues can be found in Table 1. For each code  
 250 smell, in Table 1 we also provide:

- 251 • the classification type: whether the reported code  
 252 smell is just a *suggestion*, where the choice of adopting  
 253 a correction depends on context, or it is a more serious  
 254 issue, posing a significant *problem* for the pipeline and  
 255 having a widely recognized better approach to avoid  
 256 its potential negative consequences;
- 257 • the detection method: whether the issue can be identi-  
 258 fied by using a purely *syntactic* analysis or it requires a

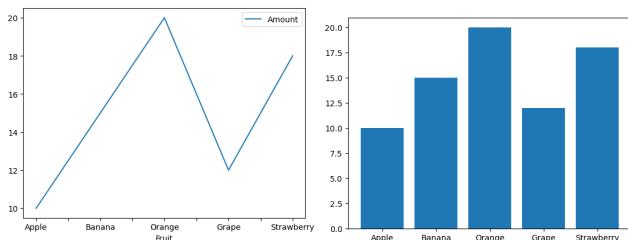
```
In [1]: import matplotlib.pyplot as plt
import pandas as pd

df = pd.read_csv("data.csv")

# DataFrame df with columns: 'Fruit', 'Amount'
# Values:
# [Apple-10, Banana-15, Orange-20,
#  Grape-12, Strawberry-18]

In [2]: # code smell: line plot
plt.plot(df["Fruit"], df["Amount"])

In [3]: # correct code
plt.bar(df["Fruit"], df["Amount"])
```



**Figure 1:** On the left, a line plot relating a string-type column and an integer-type column of a DataFrame. No exception is raised, although this plot can be deemed inadequate. On the right, a bar plot providing an appropriate visualization.

263 deeper *semantic* approach, also considering the prove-  
264 nance and content of the data;

- 265 • the **severity** level (*low, medium, high*) of the issue,  
266 based on its potential impact on the pipeline and the  
267 influence it may have on the results.

268 For clarity, we categorize the code smells into four  
269 groups: misleading visualizations, misleading results, chal-  
270 lenges for reproducibility, and general issues.

**Table 1**

Warning description (alphabetical order).

Name	Description	Type	Method	Severity Level	Severity Explanation
<b>Misleading visualizations</b>					
CategoricalPlot	A line plot is being used with categorical (nominal-scale) data on the x-axis	Suggestion	Semantic	Medium	This visualization can mislead users into interpreting categorical data as continuous, suggesting inappropriate concepts such as trends, interpolation, or monotonicity. A bar chart or similar categorical plot type should be used instead
PCAVisualization	PCA used to reduce dimensionality and visualize the data	Suggestion	Semantic	Low	PCA is not always the most appropriate technique for visualizing data
<b>Misleading results</b>					
CategoricalConversionMean	A numerical average is being calculated on categorical data that has been implicitly converted to numerical codes	Problem	Semantic	Medium	Automatic conversion of categories to numeric codes could lead to unexpected or statistically meaningless results, since the numeric codes assigned to categories do not necessarily represent a quantitative relationship between the categories themselves
DataLeakage	Information outside the training set unfairly influences a machine-learning model	Problem	Semantic	High	Data leakage may cause overestimation of performance, poor generalization, and misleading insights
DuplicatesNotDropped	Duplicated rows present in a DataFrame were not removed	Suggestion	Syntactic	Medium	Duplicates may introduce data integrity issues or bias
FixedNComponentsPCA	Principal Component Analysis (PCA) with an a priori fixed number of components	Suggestion	Syntactic	Medium	These assumptions may cause loss of important information, inefficient dimensionality reduction, and failure to identify true patterns
Gmean	The arithmetic mean is computed on ratio-based data (such as speedups), where the geometric mean would provide a more accurate measure	Problem	Semantic	Medium	Arithmetic means can be misleading or overly influenced by extreme values in this context and may result in misleading results
InappropriateMissingValues	Using summary statistics in place of the missing values	Suggestion	Syntactic	Low	This approach may distort the original data distribution, affect the correlation between variables, and introduce bias
MissingData	The DataFrame contains missing values	Suggestions	Syntactic	Medium	Missing values may cause bias, reduce the quality of the analysis, and lead to incorrect conclusions
NotShuffled	The DataFrame has not been shuffled	Suggestion	Syntactic	Low	Unshuffled data may result in biased model training and overfitting
PCAOnCategorical	PCA applied to categorical data	Suggestion	Semantic	Medium	Applying PCA to categorical data may cause suboptimal results
ScaledMean	Mean on scaled data has no direct relationship to the original data	Problem	Semantic	Medium	This may cause misleading results
<b>Challenges for reproducibility</b>					
Reproducibility	The random state is not set in train_test_split or sample function calls	Suggestion	Syntactic	Medium	This can cause reproducibility issues leading to inconsistent results
<b>General issues</b>					
HighDimensionality	A large number of features (columns) relative to the number of observations (rows)	Suggestion	Syntactic	Medium	High-dimensional data may incur the curse of dimensionality
InconsistentType	The inferred abstract type is different from the user-annotated type	Suggestion	Semantic	Low	The user annotations may not be precise
NoneRetAssignment	Assignment to a variable in the lhs where the rhs evaluation returns None	Problem	Semantic	Low	This is most likely a code smell that may result in unexpected behavior or potential runtime errors

271    **3.1. Misleading visualizations**

272    To illustrate a potential issue in data visualization, let us  
 273    consider a simple yet telling example. The pandas library  
 274    offers a variety of ways to visualize data. Ideally, users  
 275    should carefully choose the kind of plot that best fits the  
 276    nature of the data at hand. However, in practice, runtime type  
 277    checks provide little to no guidance in this respect. Consider  
 278    the code shown in Figure 1 and the generated line plot shown  
 279    below, on the left of the figure: here, a string data type (the  
 280    labels of some categorical data) on the x-axis is related to a  
 281    numeric datatype on the y-axis. Even though at first glance  
 282    this plot looks reasonable, the specific choice of a *line* plot  
 283    is questionable: a line plot hints at a continuous function  
 284    modeling the relation between domain and codomain values,  
 285    so that the user is implicitly encouraged to reason about, e.g.,  
 286    function monotonicity, local minima and maxima, or even to  
 287    approximate missing values by linear interpolation. Clearly,  
 288    all of the above makes little sense if the x-axis is representing  
 289    nominal-scale (i.e., unordered) categorical data; in such a  
 290    context, a bar chart, shown in the right hand side of Figure 1,  
 291    would have been more appropriate.

292    Another example of a code smell that can lead to mis-  
 293    leading visualizations is the use of Principal Component  
 294    Analysis (PCA), a powerful dimensionality reduction ap-  
 295    proach, for visualization purposes. In detail, PCA generates  
 296    a new set of uncorrelated features whose variance is  
 297    maximized via a linear combination of the original ones.  
 298    This new variance-based representation may not be the most  
 299    meaningful for the problem at hand and it may lead to  
 300    incorrect assumptions about the patterns within the data. An  
 301    example is shown in the left plot of Figure 2, which illus-  
 302    trates that PCA fails to produce interpretable results, thus  
 303    making highly difficult the identification of clusters within  
 304    the data. Thus, quite often PCA is not the best approach  
 305    for visualizing high-dimensional data, since its linear nature  
 306    makes it less effective at capturing more complex, non-linear  
 307    patterns in the data. In contrast, other methods such as t-  
 308    distributed stochastic neighbor embedding (t-SNE) are de-  
 309    signed to manage non-linear relationships, thus making them  
 310    particularly suitable for visualizing complex datasets [19]. In  
 311    detail, while PCA solely retains the global structures of the  
 312    data, t-SNE is able to capture local ones by preserving the  
 313    relationship between each pair of objects i.e., their similarity,  
 314    in a lower dimensional space. The latter is particularly  
 315    evident if we look at the right plot of Figure 2, which, unlike  
 316    the left one, depicts clear and identifiable clusters.

317    The two above are examples of code smells leading to  
 318    data representations being misinterpreted or confusing; the  
 319    other code smell categories focus on more insidious errors,  
 320    that in principle could go completely unnoticed.

321    **3.2. Misleading results**

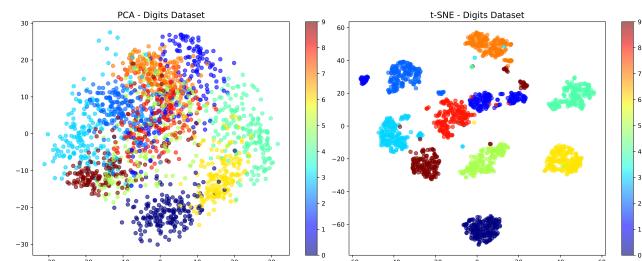
322    While being tedious for the developer, plain program-  
 323    ming errors and/or exceptions, like the one shown in Fig-  
 324    ure 3, which interrupt the normal execution flow and redirect  
 325    it to error handling code (or even program termination), are

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE

digits = datasets.load_digits()
digits_df = pd.DataFrame(data=digits.data)
digits_df['target'] = digits.target
X = digits_df.drop('target', axis=1)
y = digits_df['target']

In [2]: pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y,
            cmap='jet', alpha=0.6)

In [3]: tsne = TSNE(n_components=2,
                 perplexity=30,
                 learning_rate=200,
                 n_iter=1000,
                 random_state=42)
X_tsne = tsne.fit_transform(X)
plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=y,
            cmap='jet', alpha=0.6)
```



**Figure 2:** Comparison of PCA and t-SNE visualizations of the digits dataset. On the left, the plot resulting from PCA while on the right, the plot resulting from t-SNE. Redundant parts of the code related to plotting are omitted for clarity.

```
In [1]: import pandas as pd
x = ["Apple", "Orange", "Apple", "Apple",
      "Orange", "Apple"]
df = pd.DataFrame(x, columns=["Fruit"])
mean = df["Fruit"].mean()

Out [1]: ValueError: could not convert string to
          float: 'AppleOrangeAppleAppleOrangeApple'
```

**Figure 3:** An attempt to compute the mean of a string-type DataFrame column resulting in a ValueError exception.

actually beneficial: they force the developer to analyze and correct the issue that has arisen.

However, the highly dynamic nature and inherent flexi-  
 327    bility of Python, combined with the vast ecosystem of  
 328    libraries used in data science pipelines, can result in many  
 329    code smells or logical mistakes going unnoticed. This hap-  
 330    pens because the inaccurate action is still syntactically valid  
 331    and does not raise an exception: this behavior, often con-  
 332    sidered a feature of the language and its libraries, can lead  
 333    to unexpected results or errors.

```
In [1]: import pandas as pd
import numpy as np
from sklearn import StandardScaler,
accuracy_score, train_test_split,
LogisticRegression

df = pd.read_csv("data.csv")

X = df.iloc[:, :-1]
y = df.iloc[:, -1]

s = StandardScaler()

In [2]: # Code smell: data leakage
# Test info leaks into training
X_s = s.fit_transform(X)

X_tr, X_ts, y_tr, y_ts = train_test_split(X_s, y)

In [3]: # Corrected code
# Split before scaling
X_tr, X_ts, y_tr, y_ts = train_test_split(X, y)

X_tr = s.fit_transform(X_tr)
X_ts = s.transform(X_ts)

In [4]: m = LogisticRegression()
m.fit(X_tr, y_tr)
```

**Figure 4:** A code snippet demonstrating an approach that causes data leakage and the correct way to prevent it. The code is not executable as-is due to shortened imports for improved readability.

```
In [1]: import pandas as pd
from sklearn.decomposition import PCA

df = pd.read_csv("data.csv")
pca = PCA(n_components=3)
df_pca = pca.fit_transform(df)
print(df_pca)
```

**Figure 5:** An example of PCA with a fixed number of components.

335 to unintended consequences, where logical errors remain  
336 undetected and produce misleading results.

337 One of the most infamous and dangerous cases of mis-  
338 leading results is *data leakage*, which is exemplified in  
339 Figure 4. Data leakage occurs when information contained in  
340 the test set is inadvertently used to train the model. This can  
341 happen when some pre-processing procedures, such as data  
342 scaling, missing data imputation, over or under-sampling,  
343 etc., are performed prior to splitting the dataset into training  
344 and testing sets. The consequences of data leakage can be  
345 severe, as it can result in models with overly optimistic  
346 performances on the training set, but poor generalization,  
347 i.e., they perform poorly on unseen data, leading to incorrect  
348 predictions and potentially harmful decisions.

349 Another example of a code smell that can lead to mis-  
350 leading results is the use of PCA with a fixed number of

```
In [1]: import pandas as pd
import numpy as np

values = [25, 29, 28, 30, 27, np.nan, 150]
df = pd.DataFrame({'values': values})
# Median: 28.50, std dev: 49.92

df.fillna(df['values'].mean(), inplace=True)
# Median: 29.00, std dev: 45.57
```

**Figure 6:** An example of inappropriate missing values handling, where the mean is used to impute missing values and this leads to a different distribution of the data.

351 components (shown in Figure 5) or on categorical data.  
352 Indeed, it is common to set the number of components to 2 or  
353 3, especially if PCA is also used for visualization purposes,  
354 or to choose a number based on prior knowledge of the data,  
355 e.g., the number of classes. However, this approach can lead  
356 to overfitting, as the model may capture noise in the data  
357 rather than the underlying structure. To address this, it is  
358 essential to fine-tune this parameter, which can be achieved  
359 by objectively analyzing the results obtained with different  
360 number of components using various metrics, e.g., as the  
361 cumulative explained variance ratio of the components or  
362 the performance of a machine learning model. Similarly,  
363 applying PCA on categorical data can lead to misleading  
364 results, as it is designed for continuous data and may not  
365 capture the underlying structure of categorical data, resulting  
366 in sub-optimal performances. In such cases, it is preferable  
367 to use Multiple Correspondence Analysis (MCA), if all  
368 features are categorical, or mixed PCA, which is a technique  
369 combining MCA and PCA.

370 Moreover, several other issues can lead to misleading  
371 results, depending on the data itself or missing procedures.  
372 For example, this occurs when duplicates are not removed,  
373 the data is not randomly shuffled, or missing data is not  
374 handled correctly. In some contexts, failing to remove du-  
375 plicates can result in biased outcomes, as the model may  
376 learn from repeated instances rather than the actual data  
377 distribution. For example, a measurement that has been  
378 erroneously recorded twice by a sensor does not provide  
379 additional information but it only introduces redundancy and  
380 unbalances the dataset. Similarly, not shuffling the data can  
381 introduce bias, causing the model to learn patterns from the  
382 order of the data rather than its underlying distribution.

383 Missing data can also lead to biased results if not prop-  
384 erly addressed. Improper handling of missing values can  
385 alter the data distribution, leading to incorrect conclusions.  
386 For example, imputing missing values using summary statis-  
387 tics often introduces bias and skews the data distribution,  
388 e.g., the mean is highly sensitive to outliers, as shown in  
389 Figure 6. In such scenarios, it would be wiser to adopt more  
390 complex data imputation techniques, e.g., MissForest [37] or  
391 KNNImputer [41], to obtain more reliable estimates. Alter-  
392 natively, depending on the context and the ratio of missing  
393 data, one could remove either the affected sample or feature.

### 3.3. Challenges for Reproducibility

One of the reasons why data science pipelines are often difficult to reproduce is the lack of proper documentation and version control. This can lead to confusion and misunderstandings about the data, the analysis, and the results. For example, if the data is not properly documented, it may be difficult to understand how it was collected, what it represents, and how it was processed. On the other hand, even if the data is already provided, it may be difficult to reproduce the analysis if some preventive measures are not adopted. For example, some procedures are inherently random by default, therefore difficult to reproduce. In this case, it is important to set a random seed to ensure that the results are reproducible. This is especially important when using machine learning algorithms, as they often rely on randomness to initialize parameters or select subsets of data, i.e., when partitioning the dataset into training and testing sets. The randomness of many of these procedures is governed by a parameter called `random_state`, that works as follows. If `random_state` is set to an integer, the random number generator is seeded with that integer, ensuring that the same results are obtained each time the code is run. If `random_state` is set to `None` (the default value), the random number generator is initialized with a random seed, which means that the results will possibly be different each time the code is run.

### 3.4. General Issues

Finally, we also include some general issues that can occur in data science pipelines, related to the nature of the data or mistakes made by the developer. The eventuality of having a high dimensional dataset belongs to the first category, and it is a common issue in data science. High dimensionality is caused by the presence of a large number of features relative to a much lower number of samples in the dataset [2]. This not only makes data visualization more complex, but also leads to the curse of dimensionality, which comprises various issues caused by having too many features, ranging from an increased computational complexity to overfitting. A model that overfits accurately recognizes objects used during training, but fails to correctly characterize new, unseen objects, i.e., it is unable to generalize well. Specifically, in a high-dimensional scenario, overfitting is common since as the number of features grows, data become more sparse, making it more difficult to recognize new patterns. In other words, the number of samples required for a machine learning model to generalize well increases exponentially.

Another common issue arises from the use of `inplace` operations, which can lead to unexpected behavior and make the code difficult to understand. In-place operations modify the original data structure rather than creating a new one, therefore the return value of these operations is `None`. Nevertheless, the assignment of the return value to a variable is still possible, which can lead to confusion and unexpected behavior. Even if this is a legal assignment in Python, it is most likely not the intended behavior, and is therefore flagged as a code smell by PYRA.

## 4. PYRA's Overview

In this section we present our prototype analyzer PYRA, an Abstract Interpretation-based static analyzer for Jupyter notebooks. PYRA extends LYRA [43], a static analyzer originally developed for Python data science applications. LYRA supports input data usage analysis, so as to detect and report unused input data, and interval analysis, to infer the possible ranges of program variables.<sup>1</sup> PYRA builds upon LYRA by integrating several key features: it includes support for the analysis of non-annotated Python programs; it can handle a wider range of specific Python constructs, such as exceptions, with statements and `lambda` expressions; and it provides partial support for the libraries `pandas`, `numpy`, and `scikit-learn`, which are frequently used in data science applications. In the following we describe the architecture of PYRA, the proposed type analysis and the checkers we designed to detect the code smells discussed in Section 3.

### 4.1. Architecture

Figure 7 provides a high-level view of the architecture of PYRA: taking as input a Jupiter notebook and the confidence of checkers to be activated, PYRA produces as output an analysis report. The pipeline first converts the notebook into a Python program; in order to do this, PYRA implicitly assumes that the code cells contained in the notebook are executed in sequential order. Next, by simply visiting the Abstract Syntax Tree (AST) of the parsed Python code (i.e., the CFG generator is a subclass of the Python `ast.NodeVisitor` class) it constructs the corresponding Control-Flow Graph (CFG), i.e., a graphical and structured representation of all the paths that may be executed by the program.

Then, for each program point and each program variable, PYRA computes the corresponding abstract type information by running an Abstract Interpretation-based static analysis: this is obtained by a generic fixpoint (over-) approximation engine, parameterized with respect to the abstract domain modeling the properties of interest; the specific abstract domain we adopted for our type analysis is described in Section 4.2. Note that, before starting this static analysis phase, it is possible to enrich the input to PYRA by optionally providing the datasets on which the Jupyter notebook operates on (see the dotted line in Figure 7); this additional information, when available, can assist the static analysis in inferring more precise types for some of the variables. As an example, consider the code fragment shown in Figure 8:

When adopting a fully static approach, i.e., ignoring the contents of file `dataset.csv`, no useful type information can be derived for the data contained in `df` (and hence for the series indexed by `X` and `Y`). In contrast, if the user also provides as input the file `dataset.csv`, PYRA can infer that expression `df['X']` has a specific abstract type, e.g., `CategoricalSeries`; this additional type information can be usefully exploited by the PYRA checkers to issue an appropriate warning when later `df['X']` is used as the `x-axis` in plotting functions, as it happens in the last line of the example above.

<sup>1</sup>LYRA is publicly available at <https://github.com/caterinaurban/Lyra>.

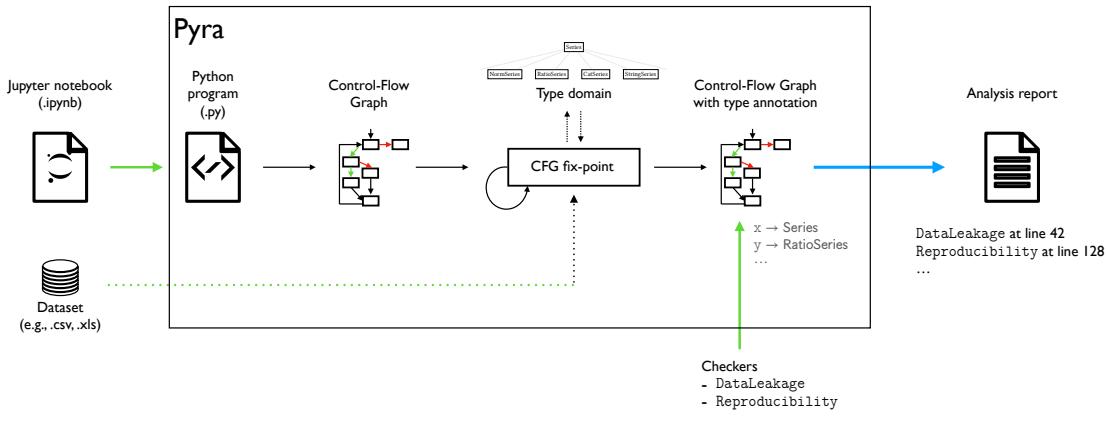


Figure 7: PYRA’s overall execution.

```
In [1]: import matplotlib.pyplot as plt
       import pandas as pd

       df = pd.read_csv("dataset.csv")
       ...
       plt.plot(df['X'], df['Y'])
```

Figure 8: Code fragment showing dataset loading and plotting.

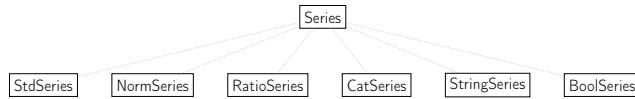


Figure 9: Diagram of the abstract domain specific to Series.

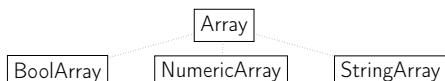


Figure 10: Diagram of the abstract domain specific to arrays.

element  $\perp$ , describing an empty set of possible values, is the most precise one and encodes a definite programming error. 519  
520

We now informally describe the elements of the abstract 521  
datatype domain used by PYRA. Currently, the domain 522  
contains 56 abstract datatypes.<sup>3</sup> 523

- Several abstract datatypes are in direct correspondence 524  
with concrete datatypes that are built-in in 525  
the language; for instance, the scalar types 526  
Bool, and String and the collection 527  
datatypes Array, List, Dict, 528  
Set, Tuple (7 abstract datatypes).
- Some special abstract datatypes for None are used 529  
in the abstract datatype domain, filtering whether 530  
None is directly assigned or is the result of an 531  
inplace operation (2 abstract datatypes).
- Other abstract datatypes are in direct correspondence 532  
with those defined in specific data science 533  
libraries, such as DataFrame and Series for 534  
pandas, or Tensor for 535  
torch.
- A few abstract datatypes are introduced to intuitively 536  
model the join of several concrete datatypes, when 537  
there seems to be no gain in keeping a fine grained 538  
differentiation; for instance, datatype Numeric 539  
is for variables storing a numeric scalar value, no matter 540  
if integral or floating point, and Scalar is for scalar 541  
values (2 abstract datatypes).
- Some abstract datatypes are introduced to model specific 542  
library functions: *encoders* (e.g., LabelEncoder, 543  
OneHotEncoder and OrdinalEncoder) are used to model 544  
scikit-learn transformers mapping the representation 545  
of categorical variables into numeric variables, so 546  
as to allow further processing (8 abstract datatypes); 547  
and *scalers*, such as StdScaler, MinMaxScaler and 548  
MaxAbsScaler (12 abstract datatypes). Consistently 549  
with our previous choices, we also model Principal 550  
Component (PCA) (1 abstract datatype).

<sup>3</sup>The full list of the PYRA’s abstract datatypes is available at [https://github.com/spangea/Pyra/blob/datascience/src/lyra/datascience/datascience\\_type\\_domain.py](https://github.com/spangea/Pyra/blob/datascience/src/lyra/datascience/datascience_type_domain.py).

<sup>2</sup>In the diagrams smaller elements are depicted below larger ones.

553      *Component Analysis (PCA)* (**1** abstract datatype),  
 554      which is used for linear dimensionality reduction by  
 555      applying a linear transformation that projects the data  
 556      into a lower-dimensional space, maximizing variance.

- 557      • Some abstract datatypes are introduced to manage  
 558      specific procedures, such as the division between the  
 559      training and test sets, which is regularly required  
 560      when developing a machine learning model (**2** abstract  
 561      datatypes). These datatypes enable our analyzer to  
 562      maintain a rather simple but sufficiently clear record  
 563      of the provenance of the data. Similarly, additional  
 564      abstract datatypes are introduced to record feature  
 565      selection, often adopted to refine the data to im-  
 566      prove performance and interpretability (**2** abstract  
 567      datatypes).
- 568      • When deemed useful, new datatypes have been in-  
 569      troduced to refine the concrete ones, so as to keep  
 570      track of relevant properties such as the way a value  
 571      has been computed. In Figure 9 we show the refine-  
 572      ments available for the `Series` datatype: for instance,  
 573      datatype `NormSeries` indicates that the values in the  
 574      series have been subjected to normalization (**8** refined  
 575      abstract datatypes for `Series`). In Figure 10 we show  
 576      the refinements for the array collections; the reason  
 577      why arrays happen to have fewer refinements with  
 578      respect to series is that they are used less frequently in  
 579      calls to the relevant data science library functions (**3**  
 580      refined abstract datatypes for `Array`). We have a similar  
 581      refinement also for list collections (**3** refined abstract  
 582      datatypes for `List`), and dataframes (**1** refined abstract  
 583      datatype for `DataFrame`).

584      In PYRA, currently, each variable is assigned a single  
 585      abstract type, although extending the analysis to a disjunc-  
 586      tive form, where each variable is mapped to a finite set of  
 587      possible types, is a possible future direction. It is also worth  
 588      highlighting that, while the current implementation of PYRA  
 589      supports 56 abstract datatypes, the framework is designed to  
 590      be easily extensible; new datatypes can be integrated into the  
 591      abstract domain by properly defining the partial order for the  
 592      newly added datatypes with respect to the already available  
 593      ones. New abstract datatypes may need to be introduced  
 594      to support the definition of new checkers, beyond those  
 595      described in the following sections.

### 4.3. Abstract Type Evaluation in PYRA

596      The static analysis computes and propagates type infor-  
 597      mation by maintaining an *abstract type environment*  $\Gamma$  that  
 598      maps each program variable  $x$  to the corresponding element  
 599       $a_x = \Gamma(x)$  of the abstract datatype domain. Intuitively,  
 600      newly encountered variables are added to  $\Gamma$  and mapped to  
 601      the top element  $T$ , meaning that nothing is initially known  
 602      about their abstract datatype; an expression  $expr$  is abstractly  
 603      evaluated to obtain its corresponding datatype, looking up  
 604      the type environment  $\Gamma$  when evaluating each of the vari-  
 605      ables occurring in the expression and combining the types

```
In [1]: import pandas as pd
        from scipy.stats import gmean
        t1 = [1.4, 5.5, 4.9, 3.9]
        t2 = [3.2, 9.8, 1.3, 1.2]

        df = pd.DataFrame({'t1': t1, 't2': t2})
        df['speedup'] = df['t1'] / df['t2']
```

Figure 11: Jupyter notebook code that shows how arithmetic mean and geometric mean can lead to different results. Since the mean is computed on speedup values, which are computed as ratios, the geometric mean is more appropriate.

of subexpressions using type rules such as

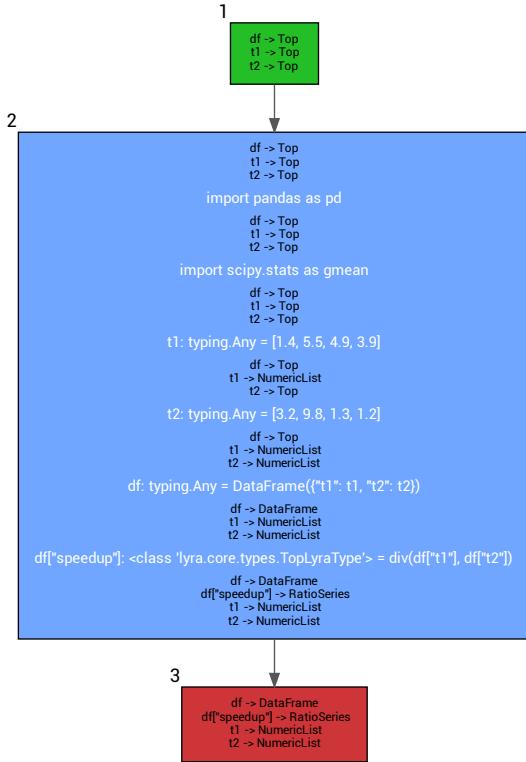
$$\text{Series} / \text{Series} = \text{RatioSeries},$$

whose intuitive reading is that the division operator, when applied to two expressions having both abstract datatype `Series`, yields a result having abstract datatype `RatioSeries`; when evaluating an assignment statement such as  $x = expr$ , we first compute the abstract datatype  $a_{expr}$  for the right-hand side expression (using  $\Gamma$ ) and then update the type environment to  $\Gamma[x \mapsto a_{expr}]$ , recording that variable  $x$  is now mapped to datatype  $a_{expr}$ . As an example, given the code fragment reported in Figure 11, PYRA produces the CFG annotated with the abstract type information shown in Figure 12; the final nodes of the CFG contain the final type information about each variable.

When joining two or more control flows, the corresponding type environments are merged by applying the abstract datatype `join` (i.e., least upper bound) operator to each variable binding; for instance, if  $\Gamma_1(x) = \text{RatioSeries}$  and  $\Gamma_2(x) = \text{StdSeries}$  then, after joining  $\Gamma_1$  and  $\Gamma_2$  into  $\Gamma$ , we obtain  $\Gamma(x) = \text{Series}$ .

*Concrete Dataset Information.* As mentioned before, it is possible to provide PYRA with the external datasets accessed and used by the Jupyter notebook. Even though not strictly necessary, this is useful to improve the precision of the analysis as it allows to compute and propagate more precise datatypes for the content of the datasets.

Algorithm 1 shows the pseudo-code of the procedure implemented in PYRA to extract abstract datatype information from the concrete dataset. The algorithm takes as input the type environment ( $\Gamma$ ), the name of the function being called (call) and the path to the dataset (path). If the call corresponds to `read_csv` (line 2), PYRA reads the CSV into a `DataFrame` using the `pandas` function (line 3) and performs several checks: whether the `DataFrame` is small (lines 4–6), high-dimensional (lines 7–9), contains duplicates (lines 10–12), or has missing values (lines 13–15), the information about these attributes is then saved (assignments at lines 5, 8, 11, 14, and 27 are kept during the analysis as further concrete information linked to the `DataFrame`) along with the abstract datatype information for the dataset in the abstract state. The values adopted for these checks are customizable and given by empirical evaluation of real-world datasets in different



**Figure 12:** PYRA’s abstract type analysis for Jupyter fragment reported in Figure 11.

**Algorithm 1** Pseudo-code of the algorithm that analyzes the concrete dataset information and maps it to the abstract datatypes.

```

1: function CONCRETE INFO( $\Gamma$ , call, path)
2:   if call = read_csv then
3:     df  $\leftarrow$  pd.read_csv(path)
4:     if LEN(df.rows)  $\leq$  100 then
5:       isSmall  $\leftarrow$  True
6:     end if
7:     if LEN(df.rows) < 2 * LEN(df.columns) then
8:       isHighDim  $\leftarrow$  True
9:     end if
10:    if HASDUPLICATES(df) then
11:      hasDuplicates  $\leftarrow$  True
12:    end if
13:    if HASNA(df) then
14:      hasNa  $\leftarrow$  True
15:    end if
16:    sortingInfo  $\leftarrow$   $\emptyset$ 
17:    for col  $\in$  df.columns do
18:      if col.dtype  $\in$  {int, float} then
19:         $\Gamma$ (col)  $\leftarrow$  NumericSeries
20:        sortingInfo[col]  $\leftarrow$  GETSORTING-
INFO(col)
21:      else if col.dtype = object then
22:         $\Gamma$ (col)  $\leftarrow$  CatSeries
23:      end if
24:    end for
25:    isShuffled  $\leftarrow$  True
26:    for col  $\in$  sortingInfo do
27:      if sortingInfo[col]  $\in$  {increasing,
decreasing} then
28:        isShuffled  $\leftarrow$  False
29:        break
30:      end if
31:    end for
32:  end if
33: end function

```

*Warning Interpretation.* In PYRA, warnings are categorized as either *plausible* or *potential* depending on the confidence of the static analysis. A *plausible* warning is emitted when the analysis has sufficient evidence to indicate that a code smell or issue is likely to occur. In contrast, a *potential* warning is issued when the analysis cannot fully determine the nature of the data or operations involved, but there are indications that a problematic pattern might be present. This distinction allows the tool to provide useful and tailored feedback, according to the desired level of confidence that can be set by the user when running PYRA.

CategoricalConversionMean, GMean, ScaledMean. Algorithm 2 reports the pseudo-code of the PYRA checker for identifying CategoricalConversionMean, GMean, and ScalerMean code smells. The checker takes as input the type environment  $\Gamma$  that occurs before the execution of the Python call. If

contexts. The algorithm also determines the datatypes of each column (lines 18–22) and assigns them to their corresponding abstract datatypes in  $\Gamma$  (lines 19 and 22). Finally, it checks if the DataFrame is shuffled based on the sorting information of its columns (lines 25–29). Note that some procedures like HASDUPLICATES (line 10) and HASNA (line 13) are omitted for brevity, but they correspond to simple checks easily implemented using the pandas library.

It is worth highlighting that providing PYRA with the dataset is not mandatory. Even if the dataset is not provided, PYRA can still analyze the code and issue warnings based on the abstract datatypes statically inferred from the code itself. Finally, independently of the dataset being provided or not, the abstract datatype for the variable related to the dataset (the left hand side of the assignment in which the right hand side is the call to `read_csv`) will always be set to `DataFrame`.

While these checks are not strictly required for the analysis to proceed, they help improve the precision and provide more information to the user about the contents of the dataset, which would otherwise remain statically unknown.

#### 4.4. PYRA Checkers

The results of the abstract type analysis are used by the checkers to identify the potential errors and code smells described in Section 3; in the following, we describe how PYRA leverages this analysis to detect them.

689 the call corresponds to `mean`, the caller `cl` is extracted (lines  
 690 2–3). Then, the abstract datatype of `cl` is retrieved from  $\Gamma$   
 691 and analyzed to generate potential warnings. Specifically,  
 692 if the abstract datatype is a `Series` datatype (line 4), then  
 693 the checker verifies whether `cl` is a `RatioSeries`, `CatSeries`,  
 694 or `ScaledSeries`. If so, a plausible related warning is issued  
 695 on that call (lines 5–9). Otherwise, the static analysis does  
 696 not have enough information to determine the exact `Series`'s  
 697 subtype of `cl`, so three potential warnings are issued (lines  
 698 12–16). Except for these cases, no warnings are raised.

---

**Algorithm 2** Pseudo-code of the `mean`'s warning-related checker.
 

---

```

1: function CHECKER( $\Gamma$ , call)
2:   if call = mean then
3:     cl  $\leftarrow$  GETCALLER(call)
4:     if  $\Gamma(cl) \subseteq \text{Series}$  then
5:       if  $\Gamma(cl) = \text{RatioSeries}$  then
6:         GMEANWARN(call, plausible)
7:       else if  $\Gamma(cl) = \text{CatSeries}$  then
8:         CATCONVMEANWARN(call, plausible)
9:       else if  $\Gamma(cl) = \text{ScaledSeries}$  then
10:        SCALEDMEANWARN(call, plausible)
11:       end if
12:     else if  $\Gamma(cl) \in \{\text{Series}, \text{T}\}$  then
13:       GMEANWARN(call, potential)
14:       CATCONVMEANWARN(call, potential)
15:       SCALEDMEANWARN(call, potential)
16:     end if
17:   end if
18: end function
```

---

699 Similarly, concerning `CategoricalConversionMean`, we apply  
 700 the same checker when inspecting the `median` call.

701 **CategoricalPlot.** When a Jupyter notebook plots something whose one of the axes is nominal-scale data, PYRA  
 702 uses Algorithm 3 to issue a warning.

---

**Algorithm 3** Pseudo-code of the `CategoricalPlot` checker.
 

---

```

1: function CHECKER( $\Gamma$ , call)
2:   if call = plot  $\wedge$  GETKIND(call)  $\notin \{\text{bar}, \text{barh}\}$  then
3:     for ax  $\in$  ARGS(call) do
4:       if  $\Gamma(ax) \in \{\text{StringList}, \text{StringArray}, \text{StringSeries}\}$  then
5:         CATPLOTWARN(call, plausible)
6:       else if  $\Gamma(ax) = \text{CatSeries}$  then
7:         CATPLOTWARN(call, plausible)
8:       else if  $\Gamma(ax) \in \{\text{Array}, \text{Series}, \text{T}\} \wedge$ 
9:          $\Gamma(ax) \notin \{\text{NumericSeries}, \text{NumericArray}\}$  then
10:          CATPLOTWARN(call, potential)
11:       end if
12:     end for
13:   end if
14: end function
```

---

When PYRA encounters a plot call which is not a bar plot, it iterates through the axis arguments (line 3) and inspects their abstract datatypes by querying  $\Gamma$ ; if the abstract datatype corresponds to `StringList`, `StringArray`, `StringSeries`, or `CatSeries`, a plausible warning is issued for the respective axis (lines 4–7). Otherwise, if  $\Gamma$  identifies the abstract datatype as either an `Array`, `Series`, or the top element (`T`), PYRA issues a potential warning.

DataLeakage. This checker is designed to identify potential data leakage issues. As previously explained, data leakage occurs when information from the test set is inadvertently used during the training phase, leading to overly optimistic performance estimates. The checker analyzes the abstract datatypes of the arguments involved in specific function calls and raises warnings if it detects potential data leakage.

Specifically the checker is activated when the functions `train_test_split`, `fit`, and `fit_transform` are called. The checker inspects the arguments of these function calls and checks for specific conditions that may indicate data leakage. The conditions checked by Algorithm 4 are the following:

- If the function call is `train_test_split` (lines 2–7), it checks if any of the arguments are of type `NormSeries`, `StdSeries`, or `CatSeries`, or if they are coming from a scaling or feature selection process (line 4). In this case, since the splitting into training and testing data sets is performed after the pre-processing, some training information may have leaked into the test set, therefore a warning is issued (line 5).
- If the function call is `fit` or `fit_transform` (lines 8–14), it checks if the fitting method is called on a test set (line 12) coming from a previous splitting operation. In this case, the warning is raised (line 11).

---

**Algorithm 4** Pseudo-code of the `DataLeakage`'s checker.
 

---

```

1: function CHECKER( $\Gamma$ , call)
2:   if call = train_test_split then
3:     for ax  $\in$  ARGS(call) do
4:       if  $\Gamma(ax) \in \{\text{NormSeries}, \text{StdSeries}, \text{CatSeries}\}$  then
5:         IS_SCALED(ax)  $\vee$  IS_FEATURE_SELECTED(ax) then
6:           DATALEAKAGEWARN(call, plausible)
7:         end if
8:       end for
9:     else if call  $\in \{\text{fit}, \text{fit\_transform}\}$  then
10:      for ax  $\in$  ARGS(call) do
11:        if IS_SPLITTED_TEST_DATA(ax) then
12:          DATALEAKAGEWARN(call, potential)
13:        end if
14:      end for
15:    end if
16: end function
```

---

DuplicatesNotDropped. The checker inspecting for this warning is syntactic, thus it does not rely on the abstract

datatype analysis described in Section 4.2. Specifically, during the abstract datatype computation, PYRA tracks whether the `drop_duplicates` method has been called on each DataFrame occurring in the Jupyter notebook. It is important to note that, this warning is always issued as *possible* warning. This is because a dataset may have been pre-processed to remove duplicates outside the notebook, without explicitly invoking methods such as `drop_duplicates` within the notebook source code, or because duplicates in some contexts may be relevant for representing the true data distribution. Consequently, when the `DuplicatesNotDropped` warning is raised, it should be interpreted as a suggestion rather than an actual error in the notebook.

751 FixedNComponentsPCA. The syntactic checker actives when  
 752 a PCA is created. Specifically, PYRA raises a warning if the  
 753 `n_components` parameter of `PCA` is assigned to a constant value,  
 754 as shown in Figure 13.

```
In [1]: // FixedNComponentsPCA warning
      pca = PCA(n_components=3)
      df_reduced = pca.fit_transform(df)
```

**Figure 13:** Example of fixed number of components in PCA.

755 As reported in Table 1, this warning should be interpreted  
 756 as a *code suggestion*. In particular, if domain knowl-  
 757 edge or prior experiments on the dataset, outside the analy-  
 758 ed notebook, suggest that a specific number of principal  
 759 components captures enough variance, setting `n_components`  
 760 may be justified. However, for improved adaptability across  
 761 different datasets, dynamically determining `n_components`,  
 762 such as by retaining a target percentage of explained vari-  
 763 ance, can be a more flexible approach.

764 HighDimensionality. The high-dimensionality checker can  
 765 be activated only if the user provides PYRA with the datasets,  
 766 allowing PYRA to extract relevant information about the  
 767 dataset applied in Algorithm 1. If the algorithm detects  
 768 high dimensionality, it raises a warning, suggesting that  
 769 feature selection, feature engineering, or dimensionality  
 770 reduction may be necessary for that dataset. Note that  
 771 there is no strict, formal definition of a high-dimensional  
 772 dataset: generally, they are loosely defined as those datasets  
 773 having far more features than samples [2]. In practice, the  
 774 high-dimensionality concept is both context- and technique-  
 775 dependent; e.g., consider the omics field, where differen-  
 776 tial expression analyses exploit all available features [29].  
 777 Hence, in PYRA we adopt a rule of thumb whereby a dataset  
 778 is considered high-dimensional when the number of features  
 779 is at least twice the number of objects. This can be seen as a  
 780 compromise that avoids raising too many warnings that are  
 781 false positives; we are aware that this threshold might be too  
 782 lax in some more classical contexts (e.g., when using a linear  
 783 regression model).

InappropriateMissingValues. PYRA may issue this warn-  
 784 ing when the code uses the `fillna` method to replace missing  
 785 values in a DataFrame with summary statistics (e.g., mean  
 786 or median). This issue becomes more concerning when the  
 787 DataFrame is small, as it can lead to misleading results. In  
 788 such cases, PYRA raises a potential warning.  
 789

InconsistentType. Python allows functions and variables  
 790 to be annotated with types, even though these annotations  
 791 are not enforced at runtime. However, if a variable is an-  
 792notated with a type, but PYRA infers an incompatible type,  
 793 the annotation is considered incorrect, and PYRA issues a  
 794 warning. Specifically, let  $x$  be a variable and  $T_x$  its user-  
 795 defined type annotation. PYRA raises a warning if  $T_x \sqcap$   
 $\Gamma(x) = \perp$ . However, no warning is issued if the inferred type  
 797 is compatible with the annotation. For example, as shown in  
 798 Figure 14:

```
In [1]: x : list = [1, 2, 3, 4]
```

**Figure 14:** Example of type annotation compatibility.

Here, PYRA infers the type of  $x$  as `NumericList`, which is  
 800 compatible with the annotated type `list`, so no warning is  
 801 generated.  
 802

MissingData. Similar to the high-dimensionality warning  
 803 checker, the missing data warning checker can be enabled if  
 804 the user provides PYRA with the datasets used. This allows  
 805 PYRA to inspect the dataset and detect any missing values  
 806 (e.g., `NaN`). If no `dropna` method is applied to the corre-  
 807 sponding DataFrame containing the dataset's information, a  
 808 warning is raised at the end of PYRA's execution.  
 809

NoneRetAssignment. Given an assignment of the form `lhs = rhs`,  
 810 if the abstract datatype static analysis infers that  $\Gamma(rhs)$   
 811 is `None`, PYRA raises a warning for the assignment. While  
 812 this operation does not inherently indicate an error or a code  
 813 smell, it may suggest a misunderstanding of the functions  
 814 or methods used in `rhs`. For example, let us consider the  
 815 following statement.  
 816

```
result = x.fillna(val, inplace=True)
```

The `fillna` method does not return a Series when the  
 818 `inplace=True` parameter is specified. As a result, assigning its  
 819 output to the variable `result` is likely unintended and could  
 820 lead to unexpected behavior in subsequent code.  
 821

NotShuffled. Similar to the `DuplicatesNotDropped` warn-  
 822 ing, the checker for `NotShuffled` is purely syntactic and  
 823 does not rely on abstract datatype analysis. During the ab-  
 824 stract datatype computation, PYRA tracks whether the `sample`  
 825 method has been called on each DataFrame in the Jupyter  
 826 notebook. As with the `DuplicatesNotDropped` warning, this  
 827 warning is always issued as a *possible* warning and should  
 828 be interpreted as a suggestion rather than an error. This is  
 829 because the dataset may have already been shuffled outside  
 830 the notebook or might be inherently random.  
 831

832 PCAOnCategorical. Algorithm 5 checks whether PCA is  
 833 applied to categorical data. When PYRA encounters a call  
 834 to transform, fit, or fit\_transform (line 2), it retrieves the  
 835 caller (line 3) and checks whether it is a PCA object (line  
 836 4). If so, it retrieves the first argument of the call (line 5) and  
 837 checks whether it is a DataFrame (line 6). If the argument is a  
 838 DataFrame, the algorithm iterates through its subscripts (line  
 839 7) (i.e. the Series belonging to it) and checks whether any of  
 840 them are categorical series (line 8). If so, a plausible warning  
 841 is issued (lines 9). Otherwise, if the analysis has not raised  
 842 a warning and has not enough information to determine the  
 843 type of the subscripts (lines 13-16), a potential warning is  
 844 issued (line 17).

**Algorithm 5** Pseudo-code of the PCAOnCategorical checker.

```

1: function CHECKER( $\Gamma$ , call)
2:   if call  $\in \{\text{transform, fit, fit\_transform}\}$  then
3:     cl  $\leftarrow \text{GETCALLER}(\text{call})$ 
4:     if  $\Gamma(\text{cl}) \sqsubseteq \text{PCA}$  then
5:       arg = GETFIRSTARG(call)
6:       if  $\Gamma(\text{arg}) \sqsubseteq \text{DataFrame}$  then
7:         for s  $\in \text{SUBSCRIPTS}(\text{arg})$  do
8:           if  $\Gamma(s) = \text{CatSeries}$  then
9:             PCAONCATWARN(call, plausible)
10:            warning_raised  $\leftarrow \text{True}$ 
11:          end if
12:        end for
13:        if  $\neg \text{warning\_raised}$  then
14:          no_warning  $\leftarrow \text{True}$ 
15:        end if
16:        if  $\neg \text{warning\_raised} \wedge \neg \text{no\_warning}$ 
then
17:          PCAONCATWARN(call, potential)
18:        end if
19:      end if
20:    end if
21:  end if
22: end function
```

845 PCAVisualization. As mentioned before, using the results  
 846 of a PCA to visualize the data is a common practice. How-  
 847 ever, this is not always the best choice, as shown in Figure 2.  
 848 In case this happens, our analyzer issues a warning following  
 849 the pseudo-code described in Algorithm 6. If the called  
 850 method is plot or scatter, the analyzer iterates through  
 851 the arguments of the call (line 3) and if the argument has  
 852 abstract datatype DataFrameFromPCA (line 4), meaning that is  
 853 a DataFrame resulting from the application of a PCA, then a  
 854 plausible warning is issued.

855 Reproducibility. If the random\_state parameter is not ex-  
 856 plicitly set when calling a method that allows for its setting,  
 857 such as the sample or train\_test\_split methods, PYRA raises  
 858 a reproducibility warning for the call.

**Algorithm 6** Pseudo-code of the PCAVisualization checker.

```

1: function CHECKER( $\Gamma$ , call)
2:   if call  $\in \{\text{plot, scatter}\}$  then
3:     for ax  $\in \text{ARGS}(\text{call})$  do
4:       if  $\Gamma(\text{ax}) = \text{DataFrameFromPCA}$  then
5:         PCAVISWARN(call, plausible)
6:       end if
7:     end for
8:   end if
9: end function
```

```

In [1]: import pandas as pd
        from sklearn StandardScaler, train_test_split
        KNeighborsClassifier, accuracy_score

        df = pd.read_csv("data.csv")
        # df.dropna(inplace=True)
        # df.drop_duplicates(inplace=True)
        # df = df.sample(frac=1, random_state=42)

        X = df.iloc[:, :-1]
        y = df.iloc[:, -1]

In [2]: sc = StandardScaler()
        X_sc = sc.fit_transform(X)

        X_tr, X_te, y_tr, y_te =
            train_test_split(X_sc, y, test_size=0.2)

In [3]: knn = KNeighborsClassifier(n_neighbors=3)
        knn.fit(X_tr, y_tr)
        y_pred = knn.predict(X_te)
        acc = accuracy_score(y_te, y_pred)
```

**Figure 15:** A code snippet containing different issues. Imports are shortened to fit the page and only refer to the library offering them, without the proper module.

**4.5. Running PYRA**

In this section, we provide a running example to illustrate how PYRA works. The example is a simple code that reads a dataset from a CSV file, splits it into training and test sets, and trains a KNeighborsClassifier model. The code is shown in Figure 15.

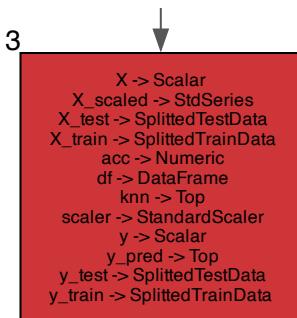
We can run PYRA on the notebook using the command:

```
pyra -analysis type-datasience code_to_analyze.py,
```

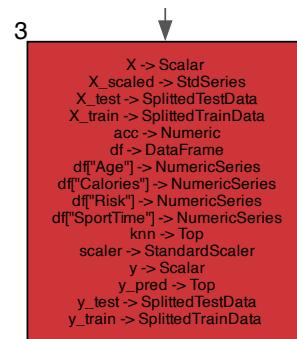
this instructs the analyzer to perform a forward analysis on the code, keeping track of the abstract datatypes and issuing both plausible and potential warnings.

The output of the analysis may vary depending on whether or not the user provides the dataset used in the code.

*Without Dataset Information.* The result of the analysis for this scenario is shown in Figure 16. In this case, PYRA is able to infer the abstract datatypes of all the variables except KNeighborsClassifier and y\_pred because their rules (i.e., the call to the constructor of KNeighborsClassifier and the call to the predict method) are not implemented in the current



**Figure 16:** PYRA’s results for the abstract type analysis of the code shown in 15 when the dataset is not provided.



**Figure 18:** PYRA’s results for the abstract type analysis of the code shown in 15 when the dataset is provided.

Reproducibility Warning			
Warning [plausible]: in train_test_split(X_sc, y, test_size=0.2) @ line 16 the random state is not set, the experiment might not be reproducible.			
Data Leakage Warning			
Warning [plausible]: in train_test_split(X_sc, y, test_size=0.2) @ line 16 data should be standardized after the split method.			

**Figure 17:** Warnings raised during the analysis of the code shown in Fig 15 when the dataset is not provided.

Age	Calories	SportTime	Risk
22	2200	4	1
<b>28</b>	<b>2100</b>	<b>Nan</b>	<b>1</b>
30	2500	5	1
<i>33</i>	<i>2400</i>	4	1
<i>33</i>	<i>2400</i>	4	1
35	2300	2	2
40	2600	2	2
<b>45</b>	<b>Nan</b>	<b>3</b>	<b>2</b>
50	2900	1	3
55	3000	0	3
60	2800	1	3

**Table 2**

Table representation of the dataset used in the running example reported in Figure 15. The rows in bold are the ones containing missing values, while the rows in italic are duplicated.

are not linked to the concrete information of the dataset. Using the information retrieved from the concrete dataset the analyzer is able to raise three new warnings. The first one is related to the presence of missing values in the dataset, and it is raised because the analyzer is able to infer that the concrete dataset contains some missing values (i.e., NaN values) and that no method has been called to drop them. The solution for this issue is to call the dropna method on the DataFrame before splitting it into training and test sets, as shown in the commented code in the snippet. The second one is related to the presence of duplicates in the dataset, which is raised because the analyzer is able to infer that the concrete dataset contains some duplicates (i.e., two rows with the same values) and that no method has been called to drop them. The solution for this issue is to call the drop\_duplicates method on the DataFrame before splitting it into training and test sets, as shown in the commented code in the snippet. Finally, the analyzer is also able to infer that the dataset is not shuffled because the first column of the dataset is sorted in increasing order. For this reason, the analyzer raises a warning suggesting to shuffle the dataset. As for the previous case, the solution for this issue is to call the sample method on the DataFrame before splitting it into training and test sets, as shown in the commented code in the snippet.

version of PYRA since they are not related to specific issues: for this reason their abstract datatypes are set to T.

Nevertheless, the analyzer is able to capture some issues and raise warnings, as shown in Figure 17. The first warning is a reproducibility issue related to the call to the `train_test_split` method without the `random_state` parameter set and it is captured with a syntactic check. This warning can be fixed by setting the `random_state` parameter to a fixed value (for example, `random_state=42`) in the arguments of the call, which is useful for reproducibility purposes. The second warning is related to a data leakage issue, which is captured by the `DataLeakage` checker (Algorithm 4). For this warning, the correct fix is similar to the one shown in Figure 4.

**With Dataset Information.** The results of the analysis when the dataset (shown in Table 2 and Figure 2) is provided are shown in Figure 18. In this case, PYRA is able to infer the abstract datatypes of all the previously detected variables that were analyzed (keeping the exception of `KNeighborsClassifier` and `y_pred`). Additionally, using the concrete analysis shown in Algorithm 1, the analyzer is able to infer the abstract datatypes of the columns of the dataset, which were not previously known, as shown in Table 2.

Moreover, based on this information and the other attributes inferred by the Algorithm 1, the analyzer is able to raise different warnings from the ones raised in the previous case, as shown in Figure 19. The issues regarding reproducibility and data leakage are still present because they

Reproducibility Warning
Warning [plausible]: in <code>train_test_split(X_sc, y, test_size=0.2)</code> @ line 16 the random state is not set, the experiment might not be reproducible.
Data Leakage Warning
Warning [plausible]: in <code>train_test_split(X_sc, y, test_size=0.2)</code> @ line 16 data should be standardized after the split method.
Missing Data Warning
Warning [potential]: At the end of the program df might still have NA values, using <code>dropna()</code> might be necessary.
Duplicates Not Dropped Warning
Warning [potential]: At the end of the program df might be small and still have duplicates that were not dropped, using <code>drop_duplicates()</code> might be necessary.
Not Shuffled Warning
Warning [potential]: At the end of the program df might be not shuffled, using <code>sample()</code> might be necessary to guarantee randomness.

**Figure 19:** Warnings raised during the analysis of the code shown in Fig 15 when the dataset is provided.

	loc	vars	calls
Minimum	21	1	6
Median	90.00	12.00	56.00
Maximum	2872	193	2123
Mean	126.84	16.45	79.58
Standard Deviation	127.33	14.71	83.32
Total	554919	71976	348181

**Table 3**  
Statistics of all the collected notebooks.

	loc	vars	calls
Minimum	21	1	6
Median	70.00	9.00	43.00
Maximum	1307	140	928
Mean	93.33	11.83	58.91
Standard Deviation	77.18	9.59	52.83
Total	204208	25875	128894

**Table 4**  
Statistics of the filtered benchmark.

Moreover, we kept only notebooks containing at least a variable, since our analyzer specifically annotates program variables, and having more than 20 lines of code (empty lines and comments are not counted), to avoid analyzing files that are too short, such as basic Kaggle templates. This criterion and these observations are meant to increase the probability that the code we analyze is somehow meaningful.

After the filtering operation, the resulting number of notebooks is 4375, and this is the benchmark on which our experimental evaluation is run. The content of these notebooks is diverse: some focus on exploratory data analysis (EDA), others build machine learning models for classification or regression tasks, while others generate visualizations or analyze patterns, and so on. The statistics on the filtered benchmark are reported in Table 4.

All the experiments were run on a 2021 MacBook Pro (model MacBookPro18,3) with M1 Pro (10 cores) and 16 GB of RAM, demonstrating how PYRA can be run on a standard laptop without requiring any special hardware or software setup. We provided the 85 projects that we used to build the benchmark, for a total of 66.76 GB of zipped data.

The processing of the entire benchmark took about 110 minutes, with an average of 2.89 seconds per notebook (minimum 1.90, maximum 106.04 seconds). This includes the time needed to analyze the notebook, as well as the time needed to unzip the folder containing the dataset and load the concrete dataset, which can be quite time consuming.

## 5.2. Qualitative Evaluation

PYRA correctly and automatically analyzes 2286 (i.e., approximately 52%) of the programs contained in the benchmark. Although this success rate may appear limited, the failures primarily arise from the intrinsic flexibility and permissiveness of Python. These features introduce challenges for static analysis tools, especially when handling highly dynamic constructs. In particular, PYRA currently

<sup>4</sup><https://www.kaggle.com/>

<sup>5</sup><https://www.kaggle.com/competitions/mayo-clinic-strip-ai>

<sup>6</sup><https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>

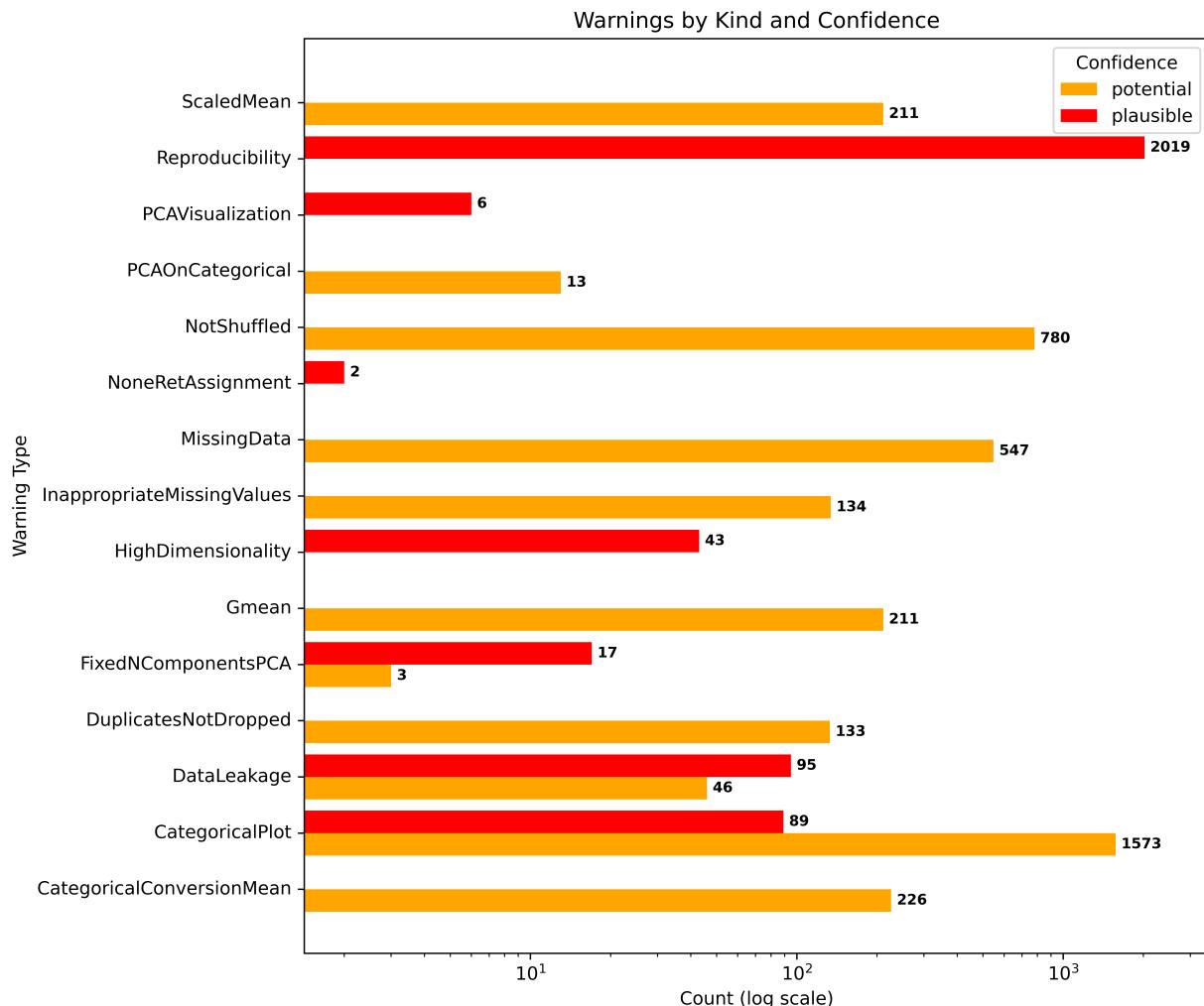


Figure 20: Warning raised in the experimental evaluation grouped by kind and confidence.

984 supports a large subset of the core language (e.g., conditional  
 985 statements, loops, exception handling), but it cannot yet  
 986 handle more intricate operations such as complex indexing  
 987 in pandas, advanced slicing mechanisms, or comprehension  
 988 constructs involving nested or dynamic expressions, which  
 989 result in exceptions. Nevertheless, it is important to highlight  
 990 that this limitation does not compromise the validity of the  
 991 proposed type analysis, being instead related to the current  
 992 prototype implementation, which still lacks support for some  
 993 advanced Python features. Further work can progressively  
 994 extend this coverage and improve the robustness of PYRA,  
 995 without requiring changes to the underlying analysis.

996 The total number of raised warnings is 4214; it is worth  
 997 noting that, even though this is a randomly collected bench-  
 998 mark, 15 of the 16 warnings that we defined were raised by  
 999 the analyzer. These warnings were found in 1661 notebooks,  
 1000 while 625 notebooks were analyzed without raising any  
 1001 warning. In detail, 50 notebooks presented warnings in 3  
 1002 out of 4 categories, while 451 had warnings in 2 of them.  
 1003 The only warning that was never raised for our benchmark

1004 is InconsistentType, only raised when the user annotates  
 1005 the type of a variable and the inferred type does not match  
 1006 the user-annotated one. Note that, type annotation is not  
 1007 a common practice in data science and its requirement is  
 1008 usually considered a constraint in the existing tools.

1009 Figure 20 shows the distribution of warnings by name  
 1010 and confidence. The most common warning was the Reproduc-  
 1011 ibility warning, which was raised 2019 times with plausible  
 1012 confidence, highlighting a significant concern regarding the  
 1013 deterministic nature of data science workflows in the ana-  
 1014 lyzed notebooks. Another of the most common warning was  
 1015 CategoricalPlot warning with a total of 1662 occurrences  
 1016 (89 plausible, 1573 potential), indicating many notebooks  
 1017 potentially misusing categorical data in plots. Related to  
 1018 the misleading visualization issue, our analysis also raised 6  
 1019 plausible PCAVisualization warnings, suggesting that some  
 1020 notebooks may not be using PCA visualizations correctly.  
 1021 Another prevalent issue was the NotShuffled warning with  
 1022 780 potential occurrences, suggesting that many data sci-  
 1023 entists may not be properly randomizing their datasets.

1024     The MissingData warning was detected 547 times with  
 1025     potential confidence, indicating notebooks that might have  
 1026     issues with missing data handling. Similarly, Categorical-  
 1027     ConversionMean warning (226 occurrences) and ScaledMean  
 1028     warning (211 occurrences) were frequently detected, both  
 1029     related to possibly improper results in statistical operations.  
 1030     The Gmean warning appeared 211 times with potential confi-  
 1031     dence.

1032     General data quality issues were also prominent, with  
 1033     DuplicatesNotDropped warning (133 occurrences) and  
 1034     InappropriateMissingValues warning (134 occurrences) sug-  
 1035     gesting that many notebooks may not properly handle data  
 1036     preprocessing steps. More critical issues like DataLeakage  
 1037     warning were detected 141 times (95 plausible, 46 potential),  
 1038     and it is worth noting that this issue could directly impact the  
 1039     performance of machine learning models.

1040     Less frequent but still significant warnings included  
 1041     HighDimensionality warning (43 occurrences),  
 1042     PCAOnCategorical warning (13 occurrences), and  
 1043     FixedNComponentsPCA warning (20 occurrences: 17 plausible,  
 1044     3 potential), all related to dimensionality or dimen-  
 1045     sionality reduction techniques. Two occurrences of the  
 1046     NoneRetAssignment warning were also detected.

1047     The wide variety and high frequency of warnings demon-  
 1048     strate the utility of PYRA in automatically detecting poten-  
 1049     tial issues in data science code that might otherwise go  
 1050     unnoticed. The distinction between potential and plausible  
 1051     warnings also provides users with information about the  
 1052     confidence level of the detected issues.

1053     It is important to emphasize that warnings with "poten-  
 1054     tial" confidence can be disabled if the user wants an analysis  
 1055     that raises less warnings. A typical use case might be when  
 1056     the user knows that certain checks are unnecessary in spe-  
 1057     cific notebooks, for example because data quality has already  
 1058     been verified in an earlier phase of the analysis or because  
 1059     some operations were intentionally performed in a certain  
 1060     way for specific purposes related to prior knowledge of the  
 1061     data. Moreover, we want to emphasize that these warnings  
 1062     are not meant to be final sentences, but rather suggestions  
 1063     for the user to consider and incentivize critical thinking  
 1064     about the code they are writing. In fact, sometimes these  
 1065     warnings need to be contextualized. For example, for the  
 1066     GMean warning it is important to take into consideration the  
 1067     distribution and scale of the data, since for logarithmic data  
 1068     the arithmetic mean might be a more appropriate choice.

### 5.2.1. Real-world Code Smells Detected by PYRA

1069     In this section, we show and discuss some examples of  
 1070     code fragments from three different notebooks contained  
 1071     in the selected benchmark suite that have raised plausible  
 1072     warnings, thus demonstrating the effectiveness of PYRA in  
 1073     identifying real-world data science code smells. The first  
 1074     one we analyze is notebook sales-eda, in which supermarket  
 1075     sales data are analyzed: first several exploratory plots are  
 1076     generated and then a Decision Tree classifier is used to  
 1077     predict customer ratings on a 1-10 scale. In Figure 21 we  
 1078     report two snippets of the notebook, that raise 6 warnings,  
 1079     4 of which are considered plausible. In detail, in the first

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
train = pd.read_csv('supermarket_sales.csv')
sns.set_theme()
plt.scatter(x = 'Branch', y = 'City',
            data = train)

In [2]: from sklearn import train_test_split
X = train_dummy.drop('Rating', axis = 1)
y = train_dummy['Rating']
X_train, X_test, y_train, y_test =
    train_test_split(X, y, test_size=0.30)
```

**Figure 21:** Example from a real notebook showing misuse of a scatter plot and reproducibility issues. Some import and names have been shortened for better readability.

```
In [1]: import pandas as pd
from sklearn import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.preprocessing import (
    MinMaxScaler, StandardScaler)
df = pd.read_csv("glass.csv")
X=df.iloc[:, :-1]
y=df.iloc[:, -1]

In [2]: minmax = MinMaxScaler()
X_minscaled = minmax.fit_transform(X)
x_minscaled
sn = []
score = []
model = DecisionTreeClassifier()
for i in range(1,101):
    X_train,X_test,y_train,y_test =
        train_test_split(X,y,stratify=y,test_size=.25)
    model.fit(X_train,y_train)
    sn.append(i)
    score.append(model.score(X_test,y_test))
```

**Figure 22:** Example from a real notebook showing reproducibil-  
 ity and data leakage issues. Some import and names have been  
 shortened for better readability.

1081 snippet, after loading data manipulation and plotting pack-  
 1082 ages, a DataFrame is created, followed by a single call to the  
 1083 scatter function from the matplotlib package. The function  
 1084 is applied to two categorical variables, Branch and City,  
 1085 making the scatter plot unsuitable: three warnings of the  
 1086 categorical plot type are raised. In the second snippet, after  
 1087 loading the necessary packages, the predictor variables x are  
 1088 defined as all columns except Rating, which is used as the  
 1089 target variable y. Then, the last line splits the original data  
 1090 into training and testing sets. However, the train\_test\_split  
 1091 function is called without setting a random seed, i.e., differ-  
 1092 ent runs can produce different partitions, thus producing a  
 1093 reproducibility issue warning.

The second notebook is classif-using-diff-scaling, which classifies different glass types using a Decision Tree model. In detail, it compares the performance of the classifier when using no standardization, z-score standardization,

```
In [1]: import pandas as pd
df_all = pd.read_csv('c19_data.csv')
df_confirmed = pd.read_csv('c19_confirmed.csv')
df_recovered = pd.read_csv('c19_recovered.csv')
df_all['datetime']=df_all['ObservationDate']
df_all['datetime']=df_all['datetime'].apply(
    lambda x:datetime.strptime(str(x),'\%m/\%d/\%Y'))
df_all['month']=df_all['datetime'].apply(
    lambda x:x.month)
df_all['day']=df_all['datetime'].apply(
    lambda x:x.day)
df_all['year']=df_all['datetime'].apply(
    lambda x:x.year)
df_all['week']=df_all['datetime'].apply(
    lambda x:x.week)
df_all['state']=df_all['Province/State']
df_all['country']=df_all['Country/Region']
df_all.drop(columns=[
    'ObservationDate','Province/State',
    'Country/Region'], inplace=True)
df_all.sample(5)
```

**Figure 23:** Example from a real notebook showing reproducibility issues. Some import and names have been shortened for better readability.

and min–max normalization. Figure 22 presents the portion of the code corresponding to the classification pipeline when employing the min–max normalization procedure. In the first code snippet, after importing the required libraries, the dataset is loaded into a DataFrame and divided into  $X$ , which contains the predictor variables, and the target variable  $y$ . The second snippet applies the min–max normalization to  $X$  and subsequently executes a loop in which the dataset is split into training and test sets, a DecisionTreeClassifier model is fit, and the corresponding accuracy is stored. Equivalent code blocks are executed for the untransformed and z-score–standardized data. Across the entire notebook, eight warnings are raised, 6 of which are classified as plausible. Two of these warnings are related to data leakage: data are normalized before being split into train and test partitions. The remaining four warnings relate to reproducibility issues caused by the random state not being set. Three of these arise from the use of the train\_test\_split function, analogously to the previous notebook, while the last one is caused by the initialization of the DecisionTreeClassifier.

The last notebook we consider is covid-19-data-analysis-and-visualization.py which presents an exploratory analysis on Covid19 data. As shown in Figure 23, it loads three CSV files into separate DataFrame objects, converts date variables into an appropriate datetime format, and extracts different date granularities, e.g., month. It also implicitly renames some columns by creating new ones and then dropping the originals. Lastly, this snippet displays the first five rows of the resulting dataset. This code actually presents 11 warnings, 3 of them plausible. Although, as mentioned, the notebook’s primary goal is exploratory, the datasets it relies on suffer from several issues, e.g., missing data, which could affect further analyses. Specifically, among the plausible warnings, two relate to high dimensional

Warning Type	Count
CategoricalPlot	6
PCAVisualization	1
CategoricalConversionMean	0
DataLeakage	16
DuplicatesNotDropped	7
FixedNComponentsPCA	2
Gmean	0
InappropriateMissingValues	7
MissingData	13
NotShuffled	16
PCAOncategorical	0
ScaledMean	0
Reproducibility	116
HighDimensionality	0
InconsistentType	0
NoneRetAssignment	0
<b>Global Statistics</b>	
Total number of warnings	184
Number of analyzed files	100
Files with warnings > 0	66
Files without warnings	34

**Table 5**

Summary of warnings and global analysis statistics.

datasets: df\_recovered and df\_confirmed are variants of the John Hopkins University CSSE COVID-19 datasets, which originally have 468 features but only 261 and 276 samples, respectively. Apart from a few location-related features, the remaining ones represent time points: comparing cities using temporal data would lead to curse of dimensionality issues. The remaining plausible issue, involves the use of the sample function without a random seed. However, in this case, the function is used just to inspect the dataset and show the newly generated fields.

### 5.3. Quantitative Evaluation

To evaluate the effectiveness of PYRA, we also randomly selected 100 notebooks from the files that PYRA correctly analyzed and manually assessed the ground truth for each file by checking the presence or absence of the issues corresponding to each warning type and cross-checking the results with all the authors. This manual assessment resulted in a total of 184 warnings across the 100 notebooks, as summarized in Table 5. The table also provides a breakdown of the number of warnings per type, along with global statistics such as the total number of warnings, the number of analyzed files, and the number of files with and without warnings. As for the qualitative analysis, also in the manual assessment, the Reproducibility warning is the most frequent one, with 116 occurrences, followed by DataLeakage (16 occurrences), showing how these two issues are particularly relevant in real-world data science code and therefore important to be detected. We then compared the warnings raised by PYRA against this ground truth to compute various performance metrics, including accuracy (Acc.), precision (Prec.), recall (Rec.), F1-score, and specificity (Spec.) for

Warning Type	Acc.	Prec.	Rec.	F1	Spec.	TP	FP	TN	FN
CategoricalConversionMean	0.941	0.000	0.000	0.000	0.941	0	6	95	0
CategoricalPlot	0.528	0.062	0.833	0.115	0.516	5	76	81	1
DataLeakage	0.922	0.833	0.625	0.714	0.977	10	2	84	6
DuplicatesNotDropped	0.970	0.833	0.714	0.769	0.989	5	1	92	2
FixedNComponentsPCA	1.000	1.000	1.000	1.000	1.000	2	0	98	0
Gmean	0.941	0.000	0.000	0.000	0.941	0	6	95	0
HighDimensionality	1.000	0.000	0.000	0.000	1.000	0	0	100	0
InappropriateMissingValues	0.970	1.000	0.571	0.727	1.000	4	0	94	3
InconsistentType	1.000	0.000	0.000	0.000	1.000	0	0	100	0
MissingData	0.950	0.722	1.000	0.839	0.943	13	5	82	0
NoneRetAssignment	1.000	0.000	0.000	0.000	1.000	0	0	100	0
NotShuffled	0.950	0.824	0.875	0.848	0.965	14	3	82	2
PCAOnCategorical	0.980	0.000	0.000	0.000	0.980	0	2	99	0
PCAVisualization	0.980	0.333	1.000	0.500	0.980	1	2	99	0
Reproducibility	0.960	0.982	0.957	0.969	0.966	111	2	56	5
ScaledMean	0.941	0.000	0.000	0.000	0.941	0	6	95	0
<b>Overall</b>	0.9256	0.5978	0.8967	0.7174	0.9290	165	111	1452	19

**Table 6**

Per-warning type metrics for combined mode (plausible + potential).

Warning Type	Acc.	Prec.	Rec.	F1	Spec.	TP	FP	TN	FN
CategoricalConversionMean	1.000	0.000	0.000	0.000	1.000	0	0	100	0
CategoricalPlot	0.922	0.000	0.000	0.000	0.979	0	2	94	6
DataLeakage	0.941	1.000	0.625	0.769	1.000	10	0	86	6
DuplicatesNotDropped	0.930	0.000	0.000	0.000	1.000	0	0	93	7
FixedNComponentsPCA	1.000	1.000	1.000	1.000	1.000	2	0	98	0
Gmean	1.000	0.000	0.000	0.000	1.000	0	0	100	0
HighDimensionality	1.000	0.000	0.000	0.000	1.000	0	0	100	0
InappropriateMissingValues	0.931	0.000	0.000	0.000	1.000	0	0	94	7
InconsistentType	1.000	0.000	0.000	0.000	1.000	0	0	100	0
MissingData	0.870	0.000	0.000	0.000	1.000	0	0	87	13
NoneRetAssignment	1.000	0.000	0.000	0.000	1.000	0	0	100	0
NotShuffled	0.842	0.000	0.000	0.000	1.000	0	0	85	16
PCAOnCategorical	1.000	0.000	0.000	0.000	1.000	0	0	100	0
PCAVisualization	0.980	0.333	1.000	0.500	0.980	1	2	99	0
Reproducibility	0.960	0.982	0.957	0.969	0.966	111	2	56	5
ScaledMean	1.000	0.000	0.000	0.000	1.000	0	0	100	0
<b>Overall</b>	0.9608	0.9538	0.6739	0.7898	0.9960	124	6	1492	60

**Table 7**

Per-warning type metrics for plausible-only mode.

both the combined levels of confidence (plausible and potential warnings) and the plausible-only level of confidence.

The overall metrics for both modes are presented in the last rows of Tables 6 and 7, respectively. These metrics are computed across all warnings raised in the 100 selected notebooks and demonstrate that PYRA performs well in both modes, with accuracy values exceeding 92%, a reasonably high F1 score exceeding 71%, and balanced precision and recall values. As expected, the plausible-only mode achieves higher precision (0.9462) but lower recall (0.6685) compared to the combined mode, which achieves a precision of

0.5942 and recall of 0.8913, reflecting the stricter criteria for raising warnings in the plausible-only mode.

A more detailed analysis is shown in Tables 6 and 7, which present the per-warning type metrics for both modes. These tables provide a detailed breakdown of the performance of PYRA for each specific warning type, allowing for a more granular analysis of its effectiveness across different types of issues.

As expected, for some warning types the results are influenced by false positives, while for others they are affected by false negatives. This is entirely anticipated, as some warnings are inherently more challenging to detect accurately through static analysis due to the complexity of the underlying issues they represent, while others may have

### MLScent vs PYRA Comparison by Warning Type

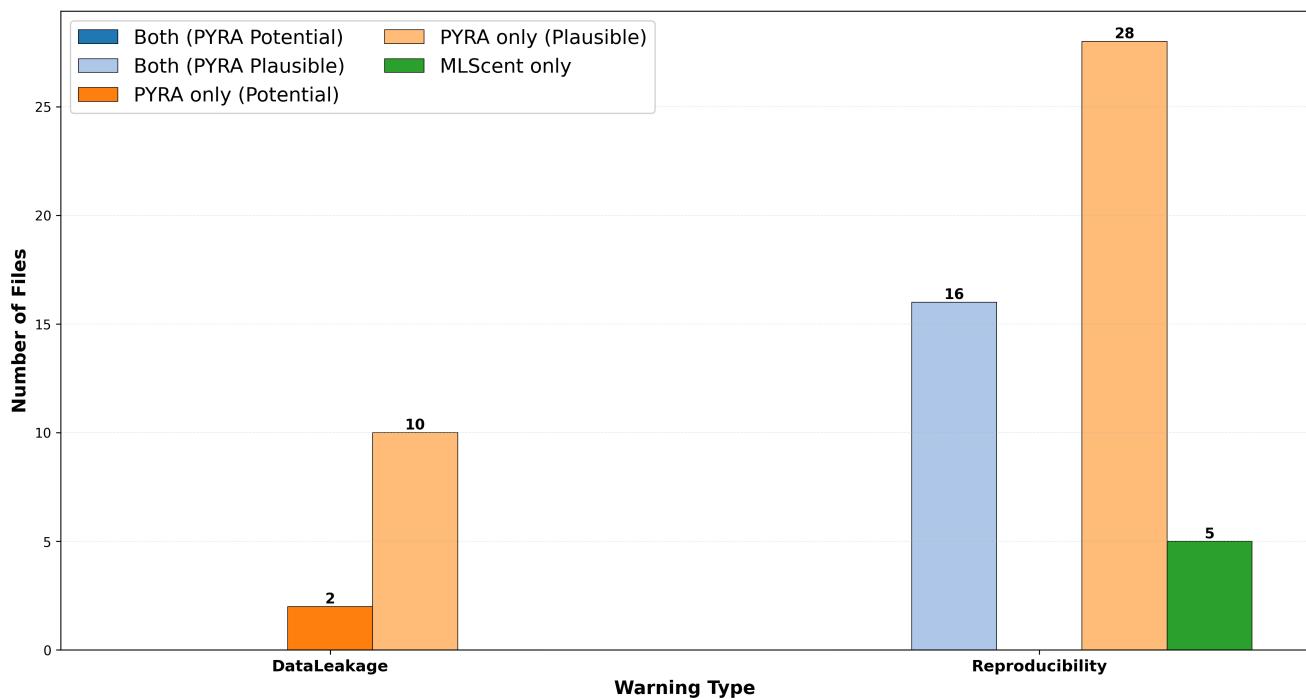


Figure 24: Comparison for DataLeakage and Reproducibility warning with MLScent.

ambiguously contexts that require user assessment for validity. For instance, the `CategoricalPlot` warning often presents difficulties in establishing a clear threshold to differentiate between correct and incorrect usage of categorical data in plots, necessitating a deep understanding of the data and analysis context, which can lead to some false positives.

Data leakage detection is also complex, with false negatives related to domain-specific knowledge (e.g., incorrect usage of time series not linked to data preprocessing) or manual operations (such as manual scaling, e.g.,  $x = (x\_data - np.min(x\_data))/(np.max(x\_data)-np.min(x\_data)).values$ ) that are not detected by static analysis. Therefore, considering the complexity of the issues being detected and the fact that some warnings have only potential confidence, the results obtained by PYRA are quite satisfactory overall, especially considering that assessing the ground truth took the authors 15 hours, while the analysis with PYRA was much faster for the entire dataset.

#### 5.4. Tool Comparison

In the quantitative evaluation benchmark, we considered the same 100 notebooks for which we manually assessed the ground truth in the quantitative evaluation and also ran another tool for detecting data science code smells, MLScent [32]. We compared its results with those of PYRA. To the best of our knowledge, there are no other publicly available tools that detect as many data science code smells as PYRA, so we focused our comparison on MLScent, which is the closest tool in terms of the number of detected code smells in common. However, the comparison can only be

made between the DataLeakage and Reproducibility warnings, as these are the only two code smells detected by both tools.

Unlike PYRA, MLScent does not provide the exact line for each warning, so we compared the results at the notebook level. Specifically, we checked whether each tool raised a warning of a given type for each notebook, regardless of the exact line where the issue was detected, and then manually validated the results.

As shown in Figure 24, PYRA outperforms MLScent in both warning types. For DataLeakage, PYRA raises this warning in 12 different files (10 with plausible confidence and 2 with potential confidence), while MLScent fails to capture any of them, even though they are all true positives. For Reproducibility, this warning is found in 16 files by both analyzers, in 28 files only by PYRA, and in 5 files only by MLScent. However, upon manually assessing these latter files, we found that they were all false positives (e.g., a warning related to a reproducibility issue for a linear regression was raised, but this operation does not involve randomness).

## 6. Discussion and Threats to Validity

Our evaluation and the design of PYRA are subject to some threats to validity. A first threat concerns false positives and false negatives. Although our experimental results show that PYRA is effective in detecting real code smells, achieving low false positive and false negative rates,

and performing favorably compared with a similar state-of-the-art tool, its precision may degrade when the dataset on which the notebook operates is not available. In such cases, PYRA falls back to a fully static approximation, reducing the precision of the inferred datatypes and potentially lowering the quality of the generated warnings. This can result in missed detections as well as spurious alerts.

Another threat arises from the assumption of sequential execution of notebook cells. While sequential execution is common and typically recommended in data-science workflows, it is not guaranteed in general. Out-of-order execution may therefore introduce discrepancies between the abstract state reconstructed by the analysis and the actual runtime behavior of the notebook.

Furthermore, PYRA currently lacks full support for some advanced Python features, such as some object-oriented programming patterns, which, although relatively uncommon in data science notebooks, may appear in more engineered workflows. As discussed in Section 5.2, this limitation does not undermine the soundness of the proposed type analysis; rather, it reflects the current state of the prototype implementation. Ongoing work is progressively extending feature coverage and improving the robustness and completeness of PYRA.

Finally, we argue that tools like the one proposed in this paper remain valuable in the era of generative AI. Indeed, such tools will be *especially* useful as data analysts increasingly rely on generative models rather than writing code themselves. We envision data analysts using PYRA to validate generated code and leveraging its analysis results and suggestions to repair the code, either manually or with the assistance of LLMs.

## 7. Conclusion

In this paper, we presented PYRA, a fully automatic static analyzer for Python data science software, aimed at detecting high-level code smells related to typical data science development pipelines rather than low-level programming errors. A key aspect of PYRA is that its warnings are designed to be easily understood not only by static analysis experts, but also, and especially, by data scientists, including early-career ones. We experimentally evaluated PYRA on a set of randomly selected real-world Jupyter notebooks crawled from Kaggle, demonstrating PYRA’s ability to detect the high-level data science issues presented and discussed in the paper, despite still being a prototype. Currently, while PYRA supports most of the core features of Python and the most popular data science libraries, some functionalities are still missing (e.g., nltk or statsmodels libraries). Future work will extend PYRA to broaden the range of Python features and libraries it supports, with the goal of increasing its applicability and usability. In this direction, we also plan to release PYRA as a plug-in for most used IDEs, such as PyCharm and Visual Studio Code.

An interesting direction for future work is to apply PYRA in the medical context, where data science plays a crucial role in tasks such as diagnosis and treatment planning. This

would involve investigating domain-specific code smells (e.g., related to data protection and privacy, or associated with the analysis of omics data) and extending PYRA with specific checkers tailored to the unique risks and code smells of medical applications. Such an extension could significantly enhance PYRA’s impact and broaden its applicability to critical, high-stakes environments.

Another promising future direction is to integrate PYRA within established quality assessment frameworks. While PYRA effectively detects code smells and potential issues, it does not by itself provide quantitative assessments of quality attributes such as maintainability, security, or reliability. Existing models for post-processing static analysis results, such as the SIG, QUAMOCO, QATCH, and SAM models [14, 25, 46, 45, 33, 34], offer mechanisms to derive actionable quality metrics. Integrating PYRA’s output within such frameworks, or developing a similar quality assessment model tailored to data science pipelines, could significantly enhance its practical value for assessing the reliability and maintainability of machine learning systems.

While we target Python, as it is currently the most popular programming language used in data science, the R programming language is also heavily used [35]. We believe that the static analyses described in this paper could be adapted to the R context as well, for instance by integrating them into flowR [36], a dataflow static analyzer for R.

Another future relevant direction could be the integration of PYRA within knowledge tracing frameworks for coding tasks, which are aimed at assessing students’ capabilities and at predicting their performances. For example in [40], large language models are used to automatically annotate knowledge concepts and PYRA could be used as an additional module to improve concept detection in Python-based data science scenarios.

Finally, at its current stage, PYRA assumes a sequential execution of notebook cells, as this is the recommended way to run a Jupyter notebook. Nevertheless, during the development phase, it is common for users to execute cells in an arbitrary order (e.g., for debugging purposes). To make PYRA applicable in such scenarios as well, a major improvement would be to support the analysis of notebooks under arbitrary execution orders.

## 8. Data Availability

The source code of PYRA is publicly available at its official Github repository: <https://github.com/spangea/Pyra>. The materials required to replicate the experimental evaluation presented in this paper are available on Zenodo at <https://zenodo.org/records/17895599>.

## Acknowledgments

This work was supported by Bando di Ateneo 2024 per la Ricerca, funded by University of Parma (FIL\_2024\_PROGETTI\_B\_IOTTI - CUP D93C24001250005).

## 1352 References

- [1] Bantilan, N., 2020. pandera: Statistical data validation of pandas dataframes, in: Agarwal, M., Calloway, C., Niederhut, D., Shupe, D. (Eds.), Proceedings of the 19th Python in Science Conference 2020 (SciPy 2020), Virtual Conference, July 6 - July 12, 2020, scipy.org. pp. 116–124. URL: <https://doi.org/10.25080/Majora-342d178e-010>, doi:10.25080/Majora-342d178e-010.
- [2] Bühlmann, P., Van De Geer, S., 2011. Statistics for high-dimensional data: methods, theory and applications. Springer Science & Business Media.
- [3] Cao, L., 2017. Data science: A comprehensive overview. ACM Comput. Surv. 50, 43:1–43:42. doi:10.1145/3076253.
- [4] Cousot, P., 1997. Types as abstract interpretations, in: Lee, P., Henglein, F., Jones, N.D. (Eds.), Conference Record of POPL'97: The 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium, Paris, France, 15-17 January 1997, ACM Press. pp. 316–331. URL: <https://doi.org/10.1145/263699.263744>, doi:10.1145/263699.263744.
- [5] Cousot, P., Cousot, R., 1977. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints, in: Graham, R.M., Harrison, M.A., Sethi, R. (Eds.), Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA, January 1977, ACM. pp. 238–252. URL: <https://doi.org/10.1145/512950.512973>, doi:10.1145/512950.512973.
- [6] Cousot, P., Cousot, R., 1992. Abstract interpretation and application to logic programs. J. Log. Program. 13, 103–179. URL: [https://doi.org/10.1016/0743-1066\(92\)90030-7](https://doi.org/10.1016/0743-1066(92)90030-7), doi:10.1016/0743-1066(92)90030-7.
- [7] Dolcetti, G., Arceri, V., Mensi, A., Zaffanella, E., Urban, C., Cortesi, A., 2025. Introducing pyra: A high-level linter for data science software, in: Dutra, I., Pechenizkiy, M., Cortez, P., Pashami, S., Pasquali, A., Moniz, N., Jorge, A.M., Soares, C., Abreu, P.H., Gama, J. (Eds.), Machine Learning and Knowledge Discovery in Databases. Applied Data Science Track and Demo Track - European Conference, ECML PKDD 2025, Porto, Portugal, September 15-19, 2025, Proceedings, Part X, Springer. pp. 449–453. doi:10.1007/978-3-032-06129-4\\_29.
- [8] Dolcetti, G., Cortesi, A., Urban, C., Zaffanella, E., 2024. Towards a high level linter for data science, in: Proceedings of the 10th ACM SIGPLAN International Workshop on Numerical and Symbolic Abstract Domains, pp. 18–25.
- [9] Drobnjakovic, F., Subotic, P., Urban, C., 2024. An abstract interpretation-based data leakage static analysis, in: Chin, W., Xu, Z. (Eds.), Theoretical Aspects of Software Engineering - 18th International Symposium, TASE 2024, Guiyang, China, July 29 - August 1, 2024, Proceedings, Springer. pp. 109–126. URL: [https://doi.org/10.1007/978-3-031-64626-3\\\_7](https://doi.org/10.1007/978-3-031-64626-3\_7), doi:10.1007/978-3-031-64626-3\\_7.
- [10] Fowler, S., Lindley, S., Morris, J.G., Decova, S., 2019. Exceptional asynchronous session types: session types without tiers. Proc. ACM Program. Lang. 3, 28:1–28:29. doi:10.1145/3290341.
- [11] Gentleman, R.C., Carey, V.J., Bates, D.M., Bolstad, B., Dettling, M., Dudoit, S., Ellis, B., Gautier, L., Ge, Y., Gentry, J., et al., 2004. Bioconductor: open software development for computational biology and bioinformatics. Genome biology 5, 1–16.
- [12] Goel, A., Donat-Bouillud, P., Krikava, F., Kirsch, C.M., Vitek, J., 2021. What we eval in the shadows: a large-scale study of eval in R programs. Proc. ACM Program. Lang. 5, 1–23. URL: <https://doi.org/10.1145/3485502>, doi:10.1145/3485502.
- [13] Hassan, M., Urban, C., Eilers, M., Müller, P., 2018. Maxsmt-based type inference for python 3, in: Chockler, H., Weissenbacher, G. (Eds.), Computer Aided Verification - 30th International Conference, CAV 2018, Held as Part of the Federated Logic Conference, FloC 2018, Oxford, UK, July 14-17, 2018, Proceedings, Part II, Springer. pp. 12–19. URL: [https://doi.org/10.1007/978-3-319-96142-2\\\_2](https://doi.org/10.1007/978-3-319-96142-2\_2), doi:10.1007/978-3-319-96142-2\\_2.
- [14] Heitlager, I., Kuipers, T., Visser, J., 2007. A practical model for measuring maintainability, in: Machado, R.J., e Abreu, F.B., da Cunha, P.R. (Eds.), Quality of Information and Communications Technology, 6th International Conference on the Quality of Information and Communications Technology, QUATIC 2007, Lisbon, Portugal, September 12-14, 2007, Proceedings, IEEE Computer Society. pp. 30–39. URL: <https://doi.org/10.1109/QUATIC.2007.8>, doi:10.1109/QUATIC.2007.8.
- [15] Kapoor, S., Narayanan, A., 2023. Leakage and the reproducibility crisis in machine-learning-based science. Patterns 4, 100804. URL: <https://doi.org/10.1016/j.patter.2023.100804>, doi:10.1016/j.PATTER.2023.100804.
- [16] Kluyver, T., et al., 2016. Jupyter notebooks – a publishing format for reproducible computational workflows, in: Loizides, F., Schmidt, B. (Eds.), Positioning and Power in Academic Publishing: Players, Agents and Agendas, IOS Press. pp. 87 – 90.
- [17] Kramm, M., Chen, R., Sudol, T., Demello, M., Caceres, A., Baum, D., Peters, A., Ludemann, P., Swartz, P., Batchelder, N., Kaptur, A., Lindzey, L., 2019. Pytype: A static type analyzer for python code. URL: <https://github.com/google/pytype>.
- [18] scikit learn.org., . Common pitfalls and recommended practices. URL: [https://scikit-learn.org/stable/common\\_pitfalls.html](https://scikit-learn.org/stable/common_pitfalls.html).
- [19] Van der Maaten, L., Hinton, G., 2008. Visualizing data using t-sne. Journal of machine learning research 9.
- [20] McKinney, W., et al., 2011. pandas: a foundational python library for data analysis and statistics. Python for high performance and scientific computing 14, 1–9.
- [21] MISRA, 2013. MISRA-C:2012 - Guidelines for the use of the C language in critical systems. MIRA Limited, Warwickshire CV10 0TU, UK.
- [22] Monat, R., Ouadjaout, A., Miné, A., 2020. Static type analysis by abstract interpretation of python programs (artifact). Dagstuhl Artifacts Ser. 6, 11:1–11:6. URL: <https://doi.org/10.4230/DARTS.6.2.11>, doi:10.4230/DARTS.6.2.11.
- [23] de Moura, L.M., Björner, N.S., 2008. Z3: an efficient SMT solver, in: Ramakrishnan, C.R., Rehof, J. (Eds.), Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings, Springer. pp. 337–340. URL: [https://doi.org/10.1007/978-3-540-78800-3\\_24](https://doi.org/10.1007/978-3-540-78800-3_24), doi:10.1007/978-3-540-78800-3\_24.
- [24] Negriani, L., Shabadi, G., Urban, C., 2023. Static analysis of data transformations in jupyter notebooks, in: Ferrara, P., Hadarean, L. (Eds.), Proceedings of the 12th ACM SIGPLAN International Workshop on the State Of the Art in Program Analysis, SOAP 2023, Orlando, FL, USA, 17 June 2023, ACM. pp. 8–13. URL: <https://doi.org/10.1145/3589250.3596145>, doi:10.1145/3589250.3596145.
- [25] Nugroho, A., Visser, J., Kuipers, T., 2011. An empirical model of technical debt and interest, in: Ozkaya, I., Kruchten, P., Nord, R.L., Brown, N. (Eds.), Proceedings of the 2nd Workshop on Managing Technical Debt, MTD 2011, Waikiki, Honolulu, HI, USA, May 23, 2011, ACM. pp. 1–8. URL: <https://doi.org/10.1145/1985362.1985364>, doi:10.1145/1985362.1985364.
- [26] Paiva, T., Damasceno, A., Figueiredo, E., Sant'Anna, C., 2017. On the evaluation of code smells and detection tools. J. Softw. Eng. Res. Dev. 5, 7. URL: <https://doi.org/10.1186/s40411-017-0041-1>, doi:10.1186/s40411-017-0041-1.
- [27] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E., 2011. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research 12, 2825–2830.
- [28] Quaranta, L., Calefato, F., Lanubile, F., 2022. Pynblint: a static analyzer for python jupyter notebooks, in: Crnkovic, I. (Ed.), Proceedings of the 1st International Conference on AI Engineering: Software Engineering for AI, CAIN 2022, Pittsburgh, Pennsylvania, May 16–17, 2022, ACM. pp. 48–49. URL: <https://doi.org/10.1145/3522664.3528612>, doi:10.1145/3522664.3528612.

- 1486 [29] Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W.,  
1487 Smyth, G.K., 2015. limma powers differential expression analyses  
1488 for rna-sequencing and microarray studies. *Nucleic acids research*  
1489 43, e47–e47.
- 1490 [30] van Rossum, G., Lehtosalo, J., Langa, L., 2014. Pep 484 – type hints.  
1491 URL: <https://peps.python.org/pep-0484/>.
- 1492 [31] Saravanan, N., Sathish, G., Balajee, J.M., 2018. Data wrangling and  
1493 data leakage in machine learning for healthcare. *JETIR- International  
1494 Journal of Emerging Technologies and Innovative Research* 5, 553–  
1495 557.
- 1496 [32] Shivashankar, K., Martini, A., 2025. Mlscent: A tool for anti-  
1497 pattern detection in ML projects, in: 4th IEEE/ACM International  
1498 Conference on AI Engineering - Software Engineering for AI, CAIN  
1499 2025, Ottawa, ON, Canada, April 27-28, 2025, IEEE. pp. 150–160.  
1500 doi:[10.1109/CAIN66642.2025.00026](https://doi.org/10.1109/CAIN66642.2025.00026).
- 1501 [33] Siavvas, M.G., Chatzidimitriou, K.C., Symeonidis, A.L., 2017.  
1502 QATCH - an adaptive framework for software product quality assess-  
1503 ment. *Expert Syst. Appl.* 86, 350–366. URL: <https://doi.org/10.1016/j.eswa.2017.05.060>,  
1504 doi:[10.1016/j.eswa.2017.05.060](https://doi.org/10.1016/j.eswa.2017.05.060).
- 1505 [34] Siavvas, M.G., Kehagias, D.D., Tzovaras, D., Gelenbe, E., 2021.  
1506 A hierarchical model for quantifying software security based on  
1507 static analysis alerts and software metrics. *Softw. Qual. J.* 29, 431–  
1508 507. URL: <https://doi.org/10.1007/s11219-021-09555-0>, doi:[10.1007/s11219-021-09555-0](https://doi.org/10.1007/s11219-021-09555-0).
- 1509 [35] Sihler, F., Pietzschmann, L., Straub, R., Tichy, M., Diera, A., Dahou,  
1510 A.H., 2025. On the anatomy of real-world R code for static analysis,  
1511 in: Koziolek, A., Lamprecht, A., Thüm, T., Burger, E. (Eds.), *Software  
1512 Engineering 2025*, Fachtagung des GI-Fachbereichs Softwaretechnik,  
1513 Karlsruhe, Germany, February 24-28, 2025, Gesellschaft für  
1514 Informatik e.V.. p. 27. URL: <https://doi.org/10.18420/se2025-27>,  
1515 doi:[10.18420/se2025-27](https://doi.org/10.18420/se2025-27).
- 1516 [36] Sihler, F., Tichy, M., 2024. flowr: A static program slicer for R,  
1517 in: Filkov, V., Ray, B., Zhou, M. (Eds.), *Proceedings of the 39th  
1518 IEEE/ACM International Conference on Automated Software Engi-  
1519 neering, ASE 2024*, Sacramento, CA, USA, October 27 - November  
1520 1, 2024, ACM. pp. 2390–2393. URL: <https://doi.org/10.1145/3691620.3695359>,  
1521 doi:[10.1145/3691620.3695359](https://doi.org/10.1145/3691620.3695359).
- 1522 [37] Stekhoven, D.J., Bühlmann, P., 2012. Missforest—non-parametric  
1523 missing value imputation for mixed-type data. *Bioinformatics* 28,  
1524 112–118.
- 1525 [38] Subotic, P., Bojanic, U., Stojic, M., 2022a. Statically detecting data  
1526 leakages in data science code, in: Gonnord, L., Titolo, L. (Eds.),  
1527 *SOAP '22: 11th ACM SIGPLAN International Workshop on the State  
1528 Of the Art in Program Analysis*, San Diego, CA, USA, 14 June 2022,  
1529 ACM. pp. 16–22. URL: <https://doi.org/10.1145/3520313.3534657>,  
1530 doi:[10.1145/3520313.3534657](https://doi.org/10.1145/3520313.3534657).
- 1531 [39] Subotic, P., Milikic, L., Stojic, M., 2022b. A static analysis framework  
1532 for data science notebooks, in: 44th IEEE/ACM International Con-  
1533 ference on Software Engineering: Software Engineering in Practice,  
1534 ICSE (SEIP) 2022, Pittsburgh, PA, USA, May 22-24, 2022, IEEE. pp.  
1535 13–22. URL: <https://doi.org/10.1109/ICSE-SEIP55303.2022.9794067>,  
1536 doi:[10.1109/ICSE-SEIP55303.2022.9794067](https://doi.org/10.1109/ICSE-SEIP55303.2022.9794067).
- 1537 [40] Sun, X., Liu, Q., Zhang, K., Shen, S., Yang, L., Li, H., 2025. Harnessing  
1538 code domain insights: Enhancing programming knowledge  
1539 tracing with large language models. *Knowledge-Based Systems*  
1540 317, 113396. URL: <https://www.sciencedirect.com/science/article/pii/S0950705125004435>,  
1541 doi:<https://doi.org/10.1016/j.knosys.2025.113396>.
- 1542 [41] Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie,  
1543 T., Tibshirani, R., Botstein, D., Altman, R.B., 2001. Missing  
1544 value estimation methods for dna microarrays. *Bioinformatics* 17,  
1545 520–525. URL: <https://doi.org/10.1093/bioinformatics/17.6.520>,  
1546 doi:[10.1093/bioinformatics/17.6.520](https://doi.org/10.1093/bioinformatics/17.6.520).
- 1547 [42] Urban, C., 2020. What programs want: Automatic inference of input  
1548 data specifications. *CoRR abs/2007.10688*. URL: <https://arxiv.org/abs/2007.10688>, arXiv:2007.10688.
- 1549 [43] Urban, C., 2023. Static analysis for data scientists, in: *Challenges of  
1550 Software Verification*. Springer, pp. 77–91.
- 1551 [44] Urban, C., Müller, P., 2018. An abstract interpretation framework for  
1552 input data usage, in: Ahmed, A. (Ed.), *Programming Languages and  
1553 Systems - 27th European Symposium on Programming, ESOP 2018,  
1554 Held as Part of the European Joint Conferences on Theory and Prac-  
1555 tice of Software, ETAPS 2018, Thessaloniki, Greece, April 14-20,  
1556 2018, Proceedings*, Springer. pp. 683–710. URL: [https://doi.org/10.1007/978-3-319-89884-1\\_24](https://doi.org/10.1007/978-3-319-89884-1_24).
- 1557 [45] Wagner, S., Goeb, A., Heinemann, L., Kläs, M., Lampasona, C.,  
1558 Lochmann, K., Mayr, A., Plösch, R., Seidl, A., Streit, J., Trendow-  
1559 icz, A., 2015. Operationalised product quality models and assess-  
1560 ment: The quamoco approach. *Inf. Softw. Technol.* 62, 101–123.  
1561 URL: <https://doi.org/10.1016/j.infsof.2015.02.009>, doi:[10.1016/j.infsof.2015.02.009](https://doi.org/10.1016/j.infsof.2015.02.009).
- 1562 [46] Wagner, S., Lochmann, K., Heinemann, L., Kläs, M., Trendowicz, A.,  
1563 Plösch, R., Seidl, A., Goeb, A., Streit, J., 2012. The quamoco product  
1564 quality modelling and assessment approach, in: Glinz, M., Murphy,  
1565 G.C., Pezzè, M. (Eds.), *34th International Conference on Software  
1566 Engineering, ICSE 2012, June 2-9, 2012, Zurich, Switzerland, IEEE  
1567 Computer Society*. pp. 1133–1142. URL: <https://doi.org/10.1109/ICSE.2012.6227106>,  
1568 doi:[10.1109/ICSE.2012.6227106](https://doi.org/10.1109/ICSE.2012.6227106).
- 1569 [47] Wang, J., Li, L., Zeller, A., 2020. Better code, better sharing: on  
1570 the need of analyzing jupyter notebooks, in: Rothermel, G., Bae, D.  
1571 (Eds.), *ICSE-NIER 2020: 42nd International Conference on Software  
1572 Engineering, New Ideas and Emerging Results*, Seoul, South Korea,  
1573 27 June - 19 July, 2020, ACM. pp. 53–56. URL: <https://doi.org/10.1145/3377816.3381724>,  
1574 doi:[10.1145/3377816.3381724](https://doi.org/10.1145/3377816.3381724).
- 1575 [48] Waskom, M.L., 2021. seaborn: statistical data visualization. *Journal  
1576 of Open Source Software* 6, 3021. doi:[10.21105/joss.03021](https://doi.org/10.21105/joss.03021).
- 1576 [49] Wickham, H., 2011. ggplot2. *Wiley interdisciplinary reviews:  
1577 computational statistics* 3, 180–185.
- 1578 [50] Zhang, H., Cruz, L., van Deursen, A., 2022. Code smells for machine  
1579 learning applications, in: Crnkoovic, I. (Ed.), *Proceedings of the 1st  
1580 International Conference on AI Engineering: Software Engineering  
1581 for AI, CAIN 2022*, Pittsburgh, Pennsylvania, May 16-17, 2022,  
1582 ACM. pp. 217–228. URL: <https://doi.org/10.1145/3522664.3528620>,  
1583 doi:[10.1145/3522664.3528620](https://doi.org/10.1145/3522664.3528620).