





#### Identifying Multi-parameter Constraint Errors in Python Data Science Library API Documentation

#### Xiufeng Xu<sup>1</sup> Fuman Xie<sup>2</sup> Chenguang Zhu<sup>3</sup> Guangdong Bai<sup>2</sup> Sarfraz Khurshid<sup>3</sup> <u>Yi Li<sup>1</sup></u>

I. Nanyang Technological University

2. University of Queensland

3. University of Texas at Austin

ISSTA '25 – Trondheim, Norway

June 25, 2025

# Code-Documentation Inconsistencies in Data Science Libraries

- Confusion<sup>[1]</sup>: Can solver **newton-cg**, **sag**, and **lbfgs** work with no penalty?
- If users use the API in a biased way, it will lead to poor model training performance like underfitting



#### penalty : {'l1', 'l2', 'elasticnet', 'none'}, default='l2'

Used to specify the norm used in the penalization. The 'newton-cg', 'sag' and 'lbfgs' solvers support only l2 penalties. 'elasticnet' is only supported by the 'saga' solver. If 'none' (not supported by the liblinear solver), no regularization is applied.

#### solver : {'newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'}, default='lbfgs'

Algorithm to use in the optimization problem.

- For small datasets, 'liblinear' is a good choice, whereas 'sag' and 'saga' are faster for large ones.
- For multiclass problems, only 'newton-cg', 'sag', 'saga' and 'lbfgs' handle multinomial loss; 'liblinear' is limited to one-versus-rest schemes.
- 'newton-cg', 'lbfgs', 'sag' and 'saga' handle L2 or no penalty

#### [1] https://github.com/scikit-learn/scikit-learn/issues/19651



Supported penalties by solver:

- 'newton-cg' ['l2', 'none']
- 'lbfgs' ['l2', 'none']
- 'liblinear' ['l1', 'l2']

- 'sag'

- ['l2', 'none']
- 'saga' ['elasticnet', 'l1', 'l2', 'none']

# Code-Documentation Inconsistencies in Data Science Libraries

- API documentation and code evolve at different speeds.
- Lots of multi-parameter constraints in data science libraries.
  - No tool do a good job of checking it !!



Constraint description of trend and seasonal in class AutoReg

> deterministic: DeterministicProcess A deterministic process. If provided, trend and seasonal are ignored. A warning is raised if trend is not "n" and seasonal is not False.

#### Corresponding code snippet in class AutoReg

- 1 class AutoReg(tsa\_model.TimeSeriesModel):
- 2 def \_\_init\_\_(...):

5

- if deterministic is not None and (self.trend != "n"(or)self.seasonal):
  - warnings.warn('When using deterministic, trend must be "n"
    - and seasonal must be False.', SpecificationWarning, stacklevel=2)





# Why is it challenging?

- Parameter-rich interfaces
  - Numerous parameters with complex dependencies
- Silent constraint violations
  - Unexpected behaviors without triggering explicit exceptions
- No fixed format for API documentation
  - Ambiguous descriptions and sometimes unclear/hidden constraints
- Implicit code-doc constraint correspondence
  - Hard to locate code segments and verify specific constraints



# LLM may be a promising solution, but ...

- LLM-only multi-parameter inconsistency checker? Probably not
  - Unavoidable stochastic behaviors and hallucination
  - Unreliable code comprehension and reasoning capability
- Our proposal: combine LLM strengths in natural language understanding with symbolic reasoning



## **Overview of MPChecker**



# Phase I: Preprocessing

- Split code into analyzable function units
- Extract documentation texts
  - Following Google or Numpy style
- Equivalent Python code transformations to support dynamic symbolic execution
  - Get rid of unnecessary complications that do not affect path constraints
    - E.g., external calls, complex data structures
    - Translate exceptions to special symbols
  - Focus on visible API parameters



#### Phase II-A: Code-constraint extraction



• Example code-constraint:

(sample\_weight != None)  $\land$  (strategy = uniform)  $\rightarrow$  ERROR



#### Phase II-B: Doc-constraint extraction

Few-shot In-context Learning

Parameter list + Few-shots+ CoT

#### Document Input Prompt





 $(affinity = nearest neighbors) \rightarrow ignore(gamma)$ 



## Implicit constraints and fuzzy words



# Phase III: Inconsistency checking

- Given a doc-constraint c and a set of codeconstraints P = {p}
  - Un-satisfiability:

$$\forall p \in P, \neg (c \land p)$$

• Nonequivalence:

$$\exists p \in P, \neg(c \Leftrightarrow p)$$



# Fuzzy constraint logic

- Hallucination  $\rightarrow$  Incorrect constraints  $\rightarrow$  fuzzy results
- LLM tends to make minor mistakes
  - Look-alike names, wrong symbols (  $< \rightarrow \leq$  ), etc.
  - Unfortunately, we can't tell whether the extracted doc-constraints are correct or not
- We could mitigate the impact of hallucination on consistency checking



- Constraint in doc-constraint: (samples < features) ^ (dual = True)</pre>
- Expression in code: n\_samples >= n\_features, dual = True



- Constraint in doc-constraint: (samples < features) ^ (dual = True)</pre>
- Expression in code: n\_samples >= n\_features, dual = True



- Constraint in doc-constraint: (samples < features) ^ (dual = True)</pre>
- Expression in code: n\_samples >= n\_features, dual = True



•  $\alpha$  and  $\beta$  denotes the relative weights. Usually, we assign the same weights

$$\sigma(e_1, e_2) = \alpha * (1 - \frac{LD(p_1, p_2)}{\max(|p_1|, |p_2|)}) + \beta * (\frac{\delta_{\mathbf{b} \neq \mathbf{i}_1} \cdot \delta_{\mathbf{b} \neq \mathbf{i}_2}}{\|\delta_{\mathbf{b} \neq \mathbf{i}_1}\| \|\delta_{\mathbf{b} \neq \mathbf{i}_2}\|}) + \alpha * (1 - \frac{LD(v_1, v_2)}{\max(|v_1|, |v_2|)})$$
Expression I: samples < features
$$\mathbf{n}_{\mathbf{s}} = \mathbf{n}_{\mathbf{s}} = \mathbf{n}_{\mathbf{s}}$$

16

### FCL – Constraint similarity

- Constraint in doc-constraint: (samples < features) ^ (dual = True)</pre>
- Expression in code: n\_samples >= n\_features, dual = True

Constraint similarity can be evaluated with the following formula:

 $\rho(c, \Phi) = \begin{cases} \arg \max \sigma(e, e_i), & \text{if } c \text{ is an expression } e \\ 1 - \sigma(c', \Phi), & \text{if } c = \neg c' \\ \min\{\sigma(c_1, \Phi), \sigma(c_2, \Phi)\}, & \text{if } c = c_1 \wedge c_2 \\ \max\{\sigma(c_1, \Phi), \sigma(c_2, \Phi)\}, & \text{if } c = c_1 \vee c_2 \end{cases} \qquad \rho = \min\{0.66, 1\} = 0.66$ 

## FCL – Membership function

• Constraint similarity serves as the degree to which a given constraint is consistent with the code.

$$\mu_{\Omega}(\epsilon) = \rho(\epsilon, \Phi_{\Omega}) \cdot \epsilon[e \mapsto e_{\Phi_{\Omega}}]$$

• After replacing each expression with its closest counterpart in path, the modified constraint is then evaluated by SMT solver to check whether the constraint is consistent with the code logic

 $0.66 \cdot True = 0.33 \cdot False$ 

- To reduce false positives, we set a constraint similarity threshold of 0.85
  - >0.85, discard the result
  - <=0.85, accept the result

# Evaluation



- Dataset<sup>[1]</sup>
  - 72 Real-world constraints from 4 popular libraries
    - scikit-learn, scipy, numpy, pandas
  - Mutation-based extended dataset (3X constraints)
    - 216 constraints (126 inconsistent + 90 consistent)
- Implementation
  - GPT-4, z3 SMT Solver
- Research Questions
  - I. How accurate is MPChecker in extracting constraints from API documentation?
  - 2. How effective is MPChecker in detecting errors related to multi-parameter constraints in API documentation?
  - 3. How effective can MPChecker detect unknown inconsistency issues?



[1] https://doi.org/10.5281/zenodo.15202267

# Evaluation: RQI

• How accurate is MPChecker in extracting constraints from API documentation?

	Equivalent	Non-Equivalent		Accuracy	
	Correct	Incorrect	Missing		
	extraction	extraction	constraints		
<b>MPCHECKER w/o few-shot learning</b>	45	20	7	62.5%	
<b>MPCHECKER w/o chain-of-thought</b>	57	7	8	79.2%	
MPCHECKER	66	2	4	<b>91.7%</b>	

- Few-shot learning and CoT can help LLM in this task
- 66/72=91.7% demonstrates MPChecker is effective in extracting logical constraints from doc

# Evaluation: RQ2

• How effective is MPChecker in detecting errors related to multiparameter constraints in API documentation?

Checker		FN	TP	Recall
LLM	raw docs + code $\rightarrow$ LLM Checker	119	7	5.6%
LLM+C	doc-constrs + code $\rightarrow$ LLM Checker	74	52	41.3%
VANILLA MPCHECKER	doc-constrs + code-constrs $\rightarrow$ Z3	39	87	69.0%
Fuzzy MPChecker	doc-constrs + code-constrs + fuzzy words + FCL $\rightarrow$ Z3	9	117	<b>92.8</b> %

- LLM performs poorly as end-to-end inconsistency checker
- LLM Results: Correct conclusion + vague / incorrect reasons
- 117/126=92.8% shows MPChecker is effective in detecting inconsistent errors
- Fuzzy words and fuzzy constraint logic improve recall by 23.8%

# Evaluation: RQ3

- How effective can MPChecker detect unknown inconsistency issues?
  - Issue report confirmation rate: 11/14 = 78.6%
  - 10/11 have already been resolved. 7 fix doc / 3 fix doc+code
  - MPChecker can detect unknown API documentation errors
  - It works even in unseen libraries which highlights its generalization capability



