What can go wrong if I change this line?

Test Failure Prediction Using Natural Language and AI

Toni Mattis
 Robert Hirschfeld

Software Architecture Group

HPI

HPI, University of Potsdam, Germany

Potsdam | 14. Nov. 2024

How do we know a Program is "correct"?



Automated Testing



"Executable documentation"

Tests and Feedback

Test feedback helps catch defects early

Running many tests delays feedback

Goal: Run relevant tests first.

Flask: 442 Tests, a few minutes Berlin-based company: 279 Tests, 18h(!)^[1]

[1] Elsner et al. 2021: https://dl.acm.org/doi/abs/10.1145/3460319.3464834

How do we know which tests are relevant* to a change?





*) fail if the change introduces a defect (regression)

Main Approaches



Coverage

Prioritize tests that cover code overlapping with the change by program analysis (static) or running tests (dynamic)

Dynamic languages

Becomes outdated

History Prioritize tests that failed previously

(for similar changes)

Requires plenty of data



Natural Language Prioritize tests with vocabulary / concepts overlapping with the change

Natural Language

Programming: "explaining to other programmers (+ your future self) how to make the computer solve a problem"



Term Importance | **TF-IDF**



Term Importance | Predictive Value

Weight terms according to how well they predicted previous test failures (precision)



occurs with a change (replaces IDF in TF-IDF)

HP

def hello_world(): **Input** (Prompt) Tokenization • Ŵ world def (hello print t Generation Select next token Language Model print Vocabulary return $P(t_n | t_{n-1}, ..., t_1)$ All known tokens " **Often GPT Architecture** Generative pre-trained Next token probabilities transformer

def hello_world(): Input (Prompt)
.

Tokenization



def hello_world(): Input (Prompt)
.

Tokenization



def hello_world(): Input (Prompt)

Tokenization



def hello_world(): **Input** (Prompt) : Tokenization def hello world print " <eos> t n **End of Sequence Token** Language Model <eos> Terminates generation ••••• $P(t_n | t_{n-1}, ..., t_1)$ \n n n

Instead of generating code, let the LLM output the probability of existing code

 $p(\text{ Existing Code }) = p_1 \times ... \times p_n$ Pretend we're generating it



Instead of generating tests, let the LLM output the probability of existing tests in response to a change. Run most probable test first.



- Format a prompt and append each test, scoring its probability
- Comment out deleted code
- Retain lexical scope

```
class AClass:
    def __init__(self):
        ...
    def method(self):
        return self.value
        return self.value + 1
        return self.value + 1
```

Change

LLM Prompt

Large Language Models | Embeddings

Vectors for (textual) data so that the **proximity** of two vectors measures the **semantic similarity** of their associated data



NLP- and Al-based Test Prioritization

Order tests so that ... run first

- 1. **TF-IDF:** Tests with (important) terms shared with changed code
- 2. Topic Model: Tests concerned with a similar set of topics
- 3. LLM: Tests most likely generated to test a change
- 4. Embedding: Tests semantically related to a change

But which one is the best?

Evaluation | Metrics

- Performance: average percentage of faults detected (APFD)
- Computes the area under the curve that plots the percentage of uncovered faults so far (y-axis) over the percentage of already executed tests (x-axis)



Obtaining Test Failures | Mutation Testing



Obtaining Test Failures | Change Mutation Testing



Results | Example: Flask

Surprise: **TF-IDF** (a very simple algorithm) sometimes beats AI



Fault detection (flask)

Results | Example: Jinja

In most larger projects, an **embedding-based** test ordering is best



Results | **Discussion**

- LLMs are "too nuanced" and currently too expensive
 - Penalize tests for bad practices ("Wouldn't have written such test")
 - Overhead: 10 20 ms GPU time > execution time of average unit test
- Often **simple heuristics** (e.g. TF-IDF) are competitive
- Limitations and Next steps
 - Models only pre-trained: Fine-tune LLM and embedding models on tests
 - Synthetic data: Run the experiments on real historical data
 - Additional features (e.g., commit messages)

Natural-language based Test Prioritization



HPI

Backup Slides

Selected Publications

- T. Mattis, L. Böhme, E. Krebs, M. Rinard, R. Hirschfeld, *Faster Feedback with AI? - A Test Prioritization Study*, Programming with AI <u>(PAI/24), Lund, Sweden, Mar. 2024</u>
- 2. <u>T. Mattis, P. Rein, F. Dürsch, and R. Hirschfeld, *RTPTorrent: An Open*source Dataset for Evaluating Regression Test Prioritization, 17th International Conference on Mining Software Repositories (MSR), Seoul, South Korea, May 2020</u>
- 3. <u>T. Mattis and R. Hirschfeld, *Lightweight Lexical Test Prioritization for Immediate Feedback*, The Art, Science, and Engineering of Programming, Vol. 4, Nr. 3.</u>
- 4. <u>D. Meier, T. Mattis, and R. Hirschfeld, *Toward Exploratory Understanding of Software Using Test Suites*, Programming Experience (PX/21), Cambridge, UK, Mar. 2021</u>

Conceptual Similarity | Topic Models



Conceptual Similarity | Topic Models



Distributional hypothesis: Terms that are **present or absent together** refer to a similar concept (topic)

HP

Conceptual Similarity | Topic Models



Conceptual Similarity | Prioritizing with Topics

Changed lines

Test case



Term Importance | **Predictive Value**

Use terms that predicted past test failures



Evaluation Details

Dataset: Python projects

- Well-tested
- Tests must run reproducibly many versions into the past (~10 years)
- Only include tests that pass an (un-mutated) control run

	Commits	Tests (Param.)	Faults	LOC Changed
Flask	159	390 (442)	726	12.5
Requests	43	314 (557)	188	13.8
Jinja	68	655 (829)	420	15.2

LLM-Based experiments:

- Performed on NVidia RTX 4090 (24GB GPU Memory)
- Models: CodeLlama-7B (initial), StableCode-3B (in the paper), CodeGemma-1.1-2b (now)

Bi-term Topic Model for Code

- Traditional topic models (e.g., LDA) designed for documents
- Code has a smaller vocabulary than natural language documents, less redundancy, and a hierarchy/graph-like structure (no apparent document boundaries)



Fitting a Code Topic Model via Edge Clustering





2024-11-14

43



