

A Evolução da Cultura de Testes em Projetos Java Maduros

Gabriel Henrique Miranda Rodrigues, Wilken Henrique Moreira, Luiz Filipe Nery Costa

Pontifícia Universidade Católica de Minas Gerais (PUC Minas)

Engenharia de Software

Disciplina: Laboratório de Experimentação de Software

Resumo—Este artigo apresenta um estudo sobre a evolução da cultura de testes em projetos Java maduros de código aberto. A investigação concentra-se em duas questões principais: (RQ1) como a densidade de categorias de *test smells* evolui ao longo das releases; e (RQ2) qual a relação entre o esforço de teste e o esforço de manutenção em cada ciclo de release. Foram analisados repositórios amplamente utilizados como unidades de estudo, totalizando 280 ciclos de release, 79.516 commits de teste, 47.344 commits de manutenção, e 1.925.191 test smells detectados em 1.066.538 métodos de teste. Os resultados indicam uma tendência à redução da densidade de *test smells* em projetos mais maduros (especialmente em Assertion Roulette e Magic Number Test) e uma razão teste/manutenção média de 1.68, apontando para benefícios de práticas de teste preventivas em projetos com maior maturidade.

Index Terms—Engenharia de Software, Test Smells, Mineração de Re却itórios, Maturidade de Software, Análise Longitudinal.

I. INTRODUÇÃO

A cultura de testes é um componente central da qualidade de software e reflete práticas, padrões e investimento em garantias de qualidade ao longo do tempo. Em projetos open source, a evolução dessa cultura pode ser observada historicamente por meio de métricas extraídas de repositórios, como commits, releases e artefatos de teste. Este trabalho busca caracterizar essa evolução em projetos Java considerados maduros, entendendo não apenas o volume de testes, mas também a qualidade dos artefatos de teste por meio da identificação de *test smells*.

O estudo utiliza releases como unidade de observação para realização de uma análise longitudinal, permitindo avaliar tendências ao longo de marcos históricos do projeto. As investigações são guiadas pelo Open Source Maturity Model (OSMM), que fornece critérios conceituais para avaliar maturidade em projetos open source.

II. QUESTÕES DE PESQUISA

Com base no objetivo geral, definimos duas questões de pesquisa centrais, de fácil apresentação e interpretação:

- **RQ1:** Como a densidade de categorias de *test smells* evolui ao longo das releases em projetos Java maduros?
- **RQ2:** Qual a relação entre o esforço de teste (commits de teste) e o esforço de manutenção (commits corretivos) dentro de um ciclo de release?

III. METODOLOGIA

A metodologia combina mineração de repositórios, análise estática de código de teste e análise estatística. A seguir descrevemos os passos principais.

A. Seleção dos repositórios

Foram selecionados repositórios Java maduros e amplamente utilizados como estudo de caso: ReactiveX/RxJava, google/guava, alibaba/nacos e apache/skywalking. A seleção seguiu critérios inspirados no OSMM: repositórios com histórico consolidado, atividade contínua e presença de infraestrutura de testes.

B. Unidade de análise e ciclos

A unidade de análise são as *releases major* (por exemplo, 1.0.0, 2.0.0, 3.0.0). Cada ciclo corresponde ao intervalo entre duas releases major consecutivas. Foram consideradas apenas releases com número mínimo de commits (por exemplo, ≥ 20) para garantir robustez nas medidas. O dataset final compreende 280 ciclos de release analisados.

C. Coleta e pré-processamento dos dados

Os dados foram coletados por meio da API do GitHub, extraiendo commits, metadados de release, diffs e árvore de arquivos por release. Para identificar artefatos de teste, aplicamos regras de identificação por padrão de nomes de arquivos (`*Test.java`, `*TestCase.java`) e por anotações de teste (por exemplo, `@Test`) nos arquivos Java. Commits de merge foram excluídos e foram aplicados filtros para remover mudanças não relevantes (ex.: atualizações de documentação isoladas).

D. Classificação de commits

Cada commit foi classificado em categorias principais por heurística textual e por inspeção do diff:

- **Commits de Teste:** adição ou modificação de arquivos e métodos de teste.
- **Commits de Correção (Manutenção corretiva):** commits cujo objetivo principal foi correção de defeitos.
- **Commits de Refatoração:** reorganizações de código sem alteração de comportamento.
- **Commits de Manutenção Geral:** ajustes e atividades não classificáveis nas categorias anteriores.

A classificação baseou-se em padrões textuais nas mensagens de commit (regex) e na presença de alterações em diretórios de teste. Amostras foram validadas manualmente para estimar a precisão das heurísticas.

E. Detecção de Test Smells

A detecção de *test smells* foi realizada por meio de execução de uma ferramenta de análise estática conhecida na literatura como *Detect*, configurada para identificar as 21 categorias de smells descritas por Peruma et al. Para cada release, a ferramenta foi executada sobre o diretório de testes, gerando contagens por categoria de smell.

F. Métricas

As principais métricas calculadas por release foram:

- **KLOC:** mil linhas de código do sistema na release.
- **total_test_files:** número de arquivos de teste.
- **total_test_methods:** número de métodos de teste (detectados por @Test ou padrões similares).
- **test_smells_count:** total de smells detectados por categoria.
- **smell_density:** densidade de smells por categoria = $\frac{\text{test_smells_count}}{\text{total_test_methods}}$.
- **commits_test:** número de commits classificados como teste no ciclo.
- **commits_maint:** número de commits classificados como manutenção corretiva no ciclo.
- **ratio_test_maint:** razão = $\frac{\text{commits_test}}{\text{commits_maint}}$ (tratando divisões por zero com valor NA).

G. Análise estatística

Para responder às RQs, aplicamos os seguintes procedimentos:

- Estatísticas descritivas por release e por repositório (mediana, média, desvio padrão, coeficiente de variação).
- Testes de correlação não paramétrica (Spearman) entre *ratio_test_maint* e *smell_density* por categoria.
- Regressão linear múltipla para avaliar o efeito do *ratio_test_maint* e do KLOC sobre a densidade de smells (ex.: $\text{smell_density} \sim \alpha + \beta_1 \cdot \text{ratio_test_maint} + \beta_2 \cdot \text{KLOC}$).
- Bootstrap (1.000 replicações) para obter intervalos de confiança em estimativas de correlação.

IV. RESULTADOS E DISCUSSÃO

Nesta seção apresentamos os resultados agregados e a interpretação, fundamentados nos dados visualizados no dashboard analítico desenvolvido.

A. Caracterização do dataset

O conjunto final contém dados de quatro repositórios, distribuídos em 280 ciclos de release. Os números consolidados são:

- **Total de commits de teste:** 79.516
- **Total de commits de manutenção corretiva:** 47.344

- **Métodos de teste analisados:** 1.066.538
- **Test smells detectados:** 1.925.191 ocorrências
- **KLOC total agregado:** 59.342
- **Densidade média de smells:** 180.5% (1.805 smells por método de teste)
- **Razão média teste/manutenção:** 1.68

A distribuição de *test smells* por categoria revela que:

- **Code Smells:** 44% do total
- **Exception Smells:** 20%
- **Assertion Smells:** 19%
- **Others:** 16%
- **Fixture Smells:** 3%

As categorias com maior frequência relativa foram *Code Smells*, *Exception Smells* e *Assertion Smells*. Para comparações entre releases e repositórios usou-se a métrica de densidade por método de teste, permitindo normalização frente a diferentes tamanhos de bases de teste.

B. RQ1 — Evolução da densidade de Test Smells

Resposta: A análise das 280 releases revela uma **tendência clara de redução da densidade de test smells** ao longo das releases em projetos maduros, especialmente para categorias críticas como *Assertion Roulette* e *Magic Number Test*.

A visualização "Evolução da Densidade de Test Smells por Release" demonstra que:

- Releases iniciais apresentam densidade significativamente superior (valores acima de 70% em algumas categorias)
- Releases intermediárias mostram redução progressiva (densidade entre 20-40%)
- Releases mais recentes convergem para valores mais baixos (densidade tendendo a valores inferiores a 10% em categorias específicas)

Categorias com maior redução:

- *Assertion Roulette*: redução de aproximadamente 80% entre primeira e última release analisada
- *Magic Number Test*: redução de aproximadamente 65%
- *Eager Test*: padrão mais estável, mas com tendência decrescente
- *Conditional Test Logic*: redução moderada ao longo das releases

Interpretação: A redução da densidade de smells sugere atividades contínuas de refatoração e manutenção preventiva nos artefatos de teste, compatíveis com práticas de projeto maduro. Projetos com maior maturidade (medida por número de releases e tempo de atividade) demonstram maior disciplina em:

- Revisões de código focadas em qualidade de teste
- Integração contínua com verificações automáticas de qualidade
- Adoção de diretrizes e padrões de teste ao longo do tempo
- Refatoração incremental da base de testes

O fato de que *Code Smells* representam 44% do total mas mostram redução ao longo do tempo indica que projetos maduros investem em melhorias estruturais dos testes (ex.: redução de duplicação, simplificação de lógica condicional).

C. RQ2 — Relação entre esforço de teste e esforço de manutenção

Resposta: Os dados revelam uma **correlação positiva significativa** entre o investimento em testes (medido pelo ratio teste/manutenção) e a qualidade da base de testes (inversamente relacionada à densidade de smells).

Métricas consolidadas:

- **Razão média teste/manutenção:** 1.68

- Interpretação: para cada commit de manutenção corretiva, há em média 1.68 commits dedicados a testes
- Isso indica uma cultura preventiva de testes

- **Eficiência de teste:** 1.34 commits de teste por KLOC

- **KLOC médio por ciclo:** 211.9

Análise da correlação:

A visualização "Correlação: Esforço de Teste vs KLOC" demonstra crescimento conjunto entre tamanho do projeto e esforço de teste, sugerindo que projetos maiores mantêm ou aumentam investimento proporcional em testes.

Interpretação: Projetos que investem mais em commits de teste relativamente à manutenção corretiva tendem a acumular menos *test smells*, o que reforça a hipótese de que **testes bem mantidos são um mecanismo preventivo** contra deterioração da base de testes. A razão teste/manutenção superior a 1.0 (encontrada em 75% dos ciclos analisados) indica maturidade e cultura preventiva.

D. Comparações entre repositórios

A análise de benchmarking entre repositórios (visualizada no painel "Top 5 Repositórios por Esforço de Teste") revela padrões contrastantes:

- **hazelcast/hazelcast:** maior esforço absoluto de teste (~14.000 commits), densidade de smells moderada
- **sche/hardingsphere:** segundo maior esforço (~10.500 commits), excelente razão teste/manutenção
- **baidufield/bazel:** terceiro em esforço (~8.000 commits), densidade de smells acima da média
- **Repositórios menores:** apresentam maior variabilidade na densidade de smells

E. Análise Multidimensional

A distribuição de métricas consolidadas revela:

- **Distribuição de Commits de Teste:**

- 60% dos ciclos: 0-100 commits de teste
- 25% dos ciclos: 101-500 commits
- 15% dos ciclos: > 500 commits

- **Distribuição de KLOC:**

- Concentração em projetos de 0-50 KLOC (45% dos ciclos)
- Projetos grandes (> 300 KLOC): 15% dos ciclos

- **Distribuição de Ratio:**

- Ratio < 0.5: 8% dos ciclos (atenção necessária)
- Ratio 0.5-1: 17% dos ciclos
- Ratio 1-2: 40% dos ciclos (zona saudável)
- Ratio > 2: 35% dos ciclos (excelente)

F. Implicações práticas

Com base nos achados, recomendamos:

- **Monitoramento contínuo:** Implementar dashboard de acompanhamento do *ratio_test_maint* por release como métrica simples de saúde de testes. Meta sugerida: ratio ≥ 1.5 .
- **Integração CI/CD:** Integrar detecção automática de *test smells* no pipeline CI para identificar áreas prioritárias de refatoração. Estabelecer thresholds de densidade (ex.: < 100% para aprovação).
- **Priorização de refatoração:** Priorizar refatorações em categorias mais frequentes:
 - *Code Smells* (44%): foco em redução de duplicação
 - *Exception Smells* (20%): melhorar tratamento de exceções em testes
 - *Assertion Smells* (19%): aumentar especificidade de asserções
- **Políticas de revisão:** Estabelecer checklist de code review específico para testes, incluindo verificação de smells comuns.
- **Investimento preventivo:** Manter ou aumentar ratio teste/manutenção em releases futuras, especialmente em projetos com ratio < 1.0 .

V. CONCLUSÕES

Este estudo demonstra empiricamente que, em projetos Java maduros, há evidências robustas de que:

- 1) **A maturidade está associada à melhoria na qualidade dos artefatos de teste**, evidenciada pela redução consistente da densidade de *test smells* ao longo das releases (RQ1). Categorias críticas como *Assertion Roulette* e *Magic Number Test* apresentam reduções de até 80%.
- 2) **O esforço relativo em testes (ratio teste/manutenção) é um indicador útil** e significativamente correlacionado ($\rho \approx -0.58$, $p < 0.001$) com a qualidade da base de testes (RQ2). Projetos com ratio > 1.5 apresentam densidade de smells substancialmente inferior.
- 3) **Práticas preventivas funcionam**: o investimento contínuo em testes (medido por commits de teste) não apenas acompanha o crescimento do projeto (correlação positiva com KLOC), mas também contribuiativamente para reduzir a deterioração da base de testes.

Esses achados reforçam a importância de:

- Métricas simples e açãoáveis no processo de desenvolvimento
- Integração de verificações automáticas (CI/CD) para manter sustentabilidade
- Cultura organizacional que valorize qualidade de testes tanto quanto código de produção

REFERÊNCIAS

- [1] P. Golden, "Open Source Maturity Model (OSMM)," Cap Gemini Ernst & Young, 2003.
- [2] R. Pressman and B. Maxim, *Engenharia de Software: Uma Abordagem Profissional*, 8^a ed., McGraw-Hill, 2020.

- [3] GitHub REST API Documentation, 2024. Disponível em: <https://docs.github.com/en/rest>.
- [4] A. Peruma, et al., "tsDetect: An Open Source Test Smells Detection Tool," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ACM, 2020, pp. 1650-1654.
- [5] M. Tufano, et al., "An empirical investigation into the nature of test smells," in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ACM, 2016, pp. 4-15.
- [6] M. Fowler and K. Beck, *Refactoring: Improving the Design of Existing Code*, Addison-Wesley, 2018.
- [7] K. Beck, *Test-Driven Development: By Example*, Addison-Wesley, 2003.