

Fehlervorhersage in Java-Code mittels Machine Learning

Tobias Meier (meierto3), Yacine Mekesser (mekesyac)

4. Juli 2016

ZHAW

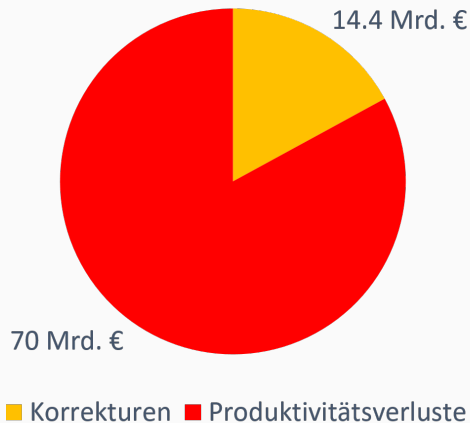
Inhaltsverzeichnis

1. Einleitung
2. Grundlagen
3. Vorgehen
4. Architektur
5. Repository Miner
6. Features
7. ML-Pipeline
8. Demo
9. Resultate
10. Zusammenfassung
11. Fragen und Diskussion

Einleitung

Softwarefehler I

- 84.4 Mrd. € jährliche Verluste in DE



Softwarefehler nehmen zu:

- Komplexität
- Wachsende Codebase
- Agile Softwareentwicklung

Frühwarnsystem für Softwarefehler

- Qualität
- Effizienz
- Günstiger

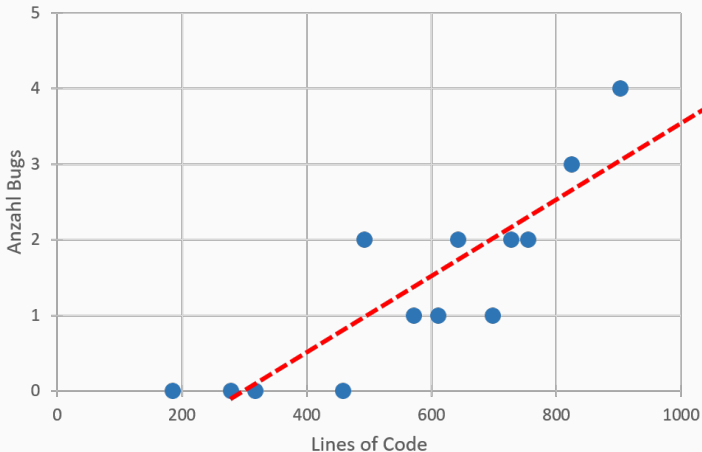
- Textverständnis mit Machine Learning
- Open-Source-Projekte
 - Öffentliche Code Repositories
 - Issue-Tracking-Systeme

- Vorhersage von Bugs
- Codeanalyse mit ML
- Beschränkung auf Java-Projekte

Grundlagen

Machine Learning

Ermöglicht Computern aus Beispielen zu **lernen** und dann zu **verallgemeinern**.



Vorgehen

Statischen Codeanalyse



checkstyle



CODE CLIMATE

Machine Learning

- Ideen und Potential vorhanden
- Keine schlüsselfertige Produkte

- Grundstein für weitere Arbeiten legen
- Toolset für projektbezogenes ML bereitstellen
- Einbeziehung von Textanalyse-Features

- ML Datensatz
 - Code und Issue-Tracking
- Features
- ML-Pipeline

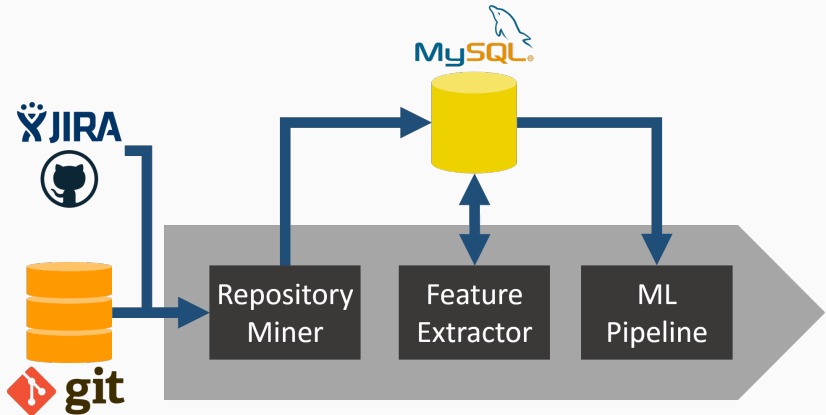
Architektur

Drei separate Tools:

- Repository Miner
- Feature Extractor
- ML-Pipeline

Ermöglicht Modularität , Flexibilität & Effizienz

Systemarchitektur

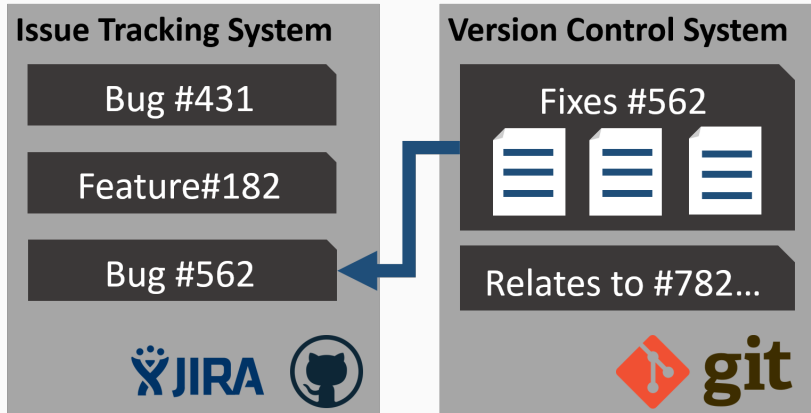


Repository Miner

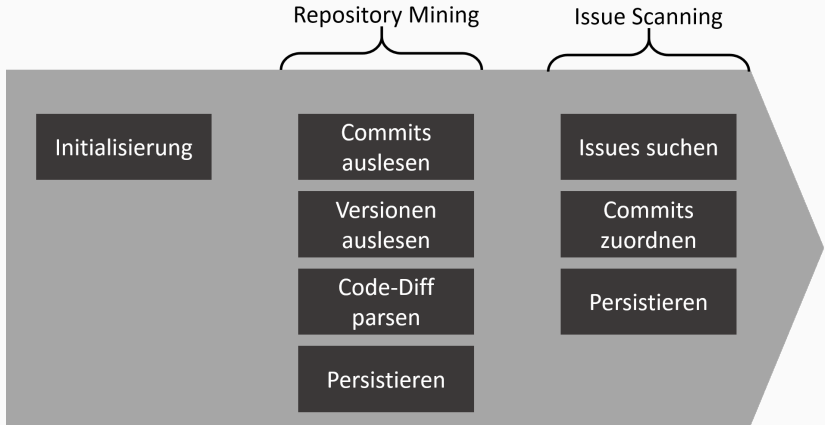
- Anbindung an Git, GitHub & Jira
- Schnellen Zugriff bieten
- Daten strukturiert und normalisiert speichern

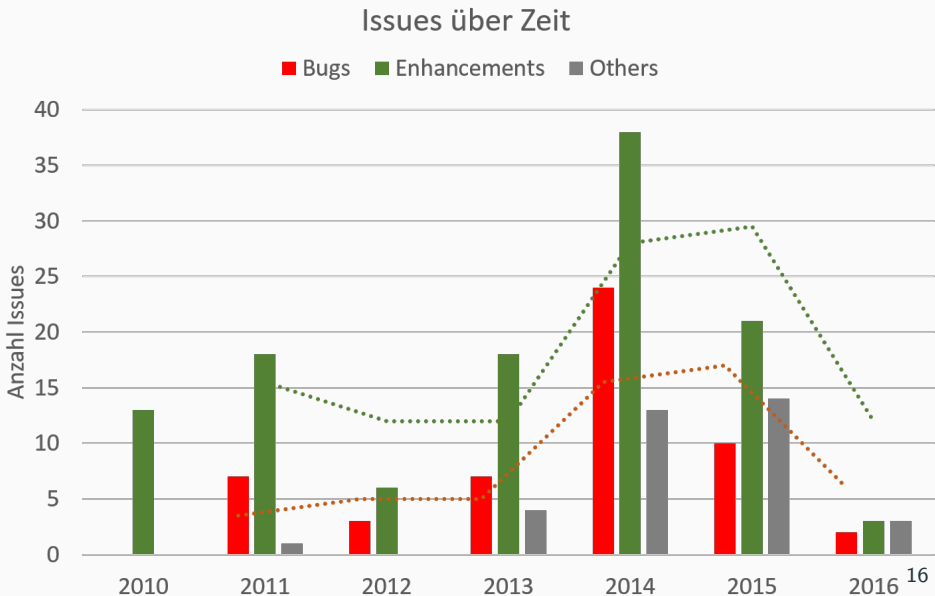


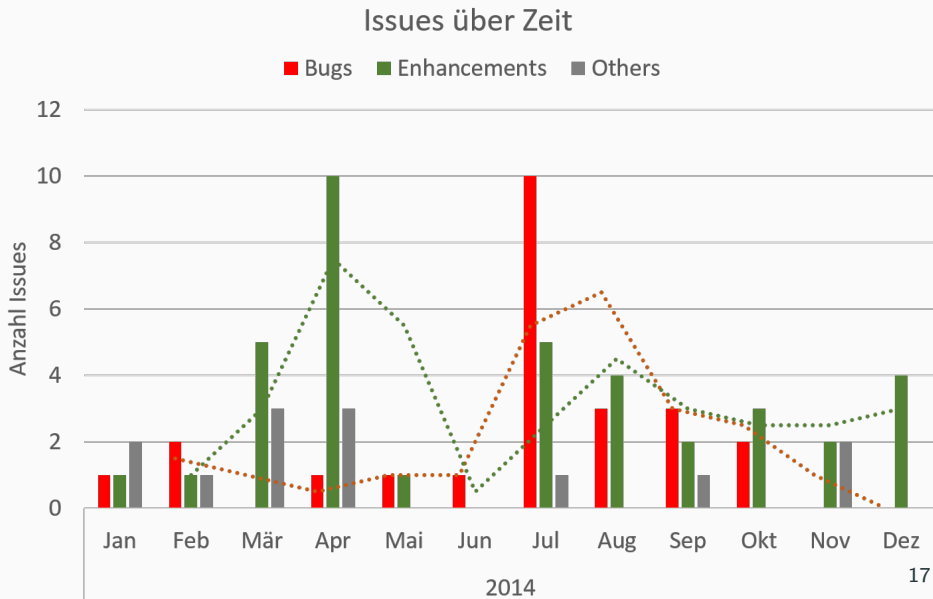
Herausforderungen



Workflow



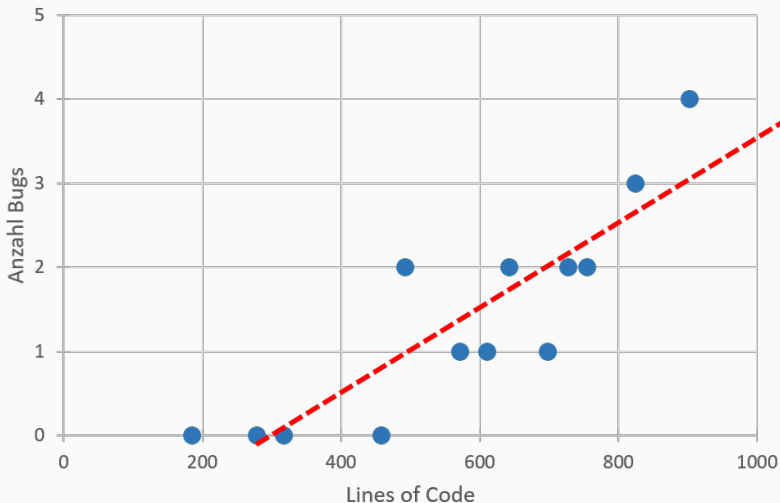




Features

Features

- Messbare Eigenschaften einer Beobachtung



Lines-of-Code-Features

- Anzahl Codezeilen
- Grosse Java Files - viele Fehler?
- Einfache Implementation

Objektorientierte Features

- Kopplung zwischen Objekten
- Mangel an Kohäsion in Klasse
- Anzahl Methoden pro Klasse

- Halstead
 - Gegenüberstellung von Operatoren & Operanden
- McCabe
 - Anzahl Kontrollfluss-Statements

Anzahl-und-Typen-Features

- Anzahl Imports, Interfaces, Methoden, Klassenvariablen
- Komplexität von Klasse

Temporale Features

- Anzahl Bugs im letzten Monat
- Anzahl Tage seit letztem Commit
- Datenbasis: Historie der Versionsverwaltung

- Länge von Namen (Methoden, Variablen)
- Text- und Sprachanalyse

N-Grams

- Bug-Mustererkennung
- Verschiedene Abstraktionslevels

```
Code (Java): public class MyClass..
```

```
Code (AST):
```

```
(55, TypeDeclaration)  
  (83, Modifier)  
    (43, SimpleType)  
      (42, SimpleName)
```

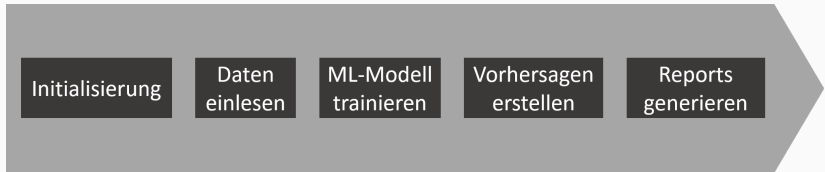
```
N-Gram: 55_83_43_42
```


- Berechnet Features aus der Datenbasis
- AST-Parsing
- Modularer Aufbau

ML-Pipeline

- Umgang mit **grossen Datensets**
 - Dataset-Caching
 - Sparse-Matrizen
- Einfache Konfiguration verschiedener **ML-Algorithmen**
 - Zentrales Config-File
- Resultate müssen **bewertet** und **analysiert** werden
 - Verschiedene Charts & Reports

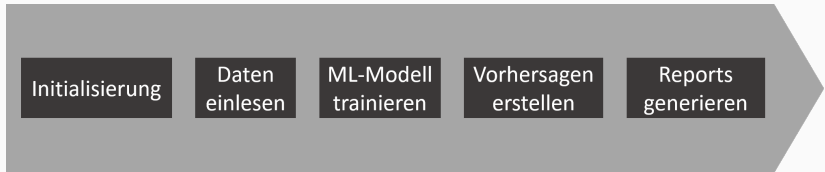
Workflow



Demo

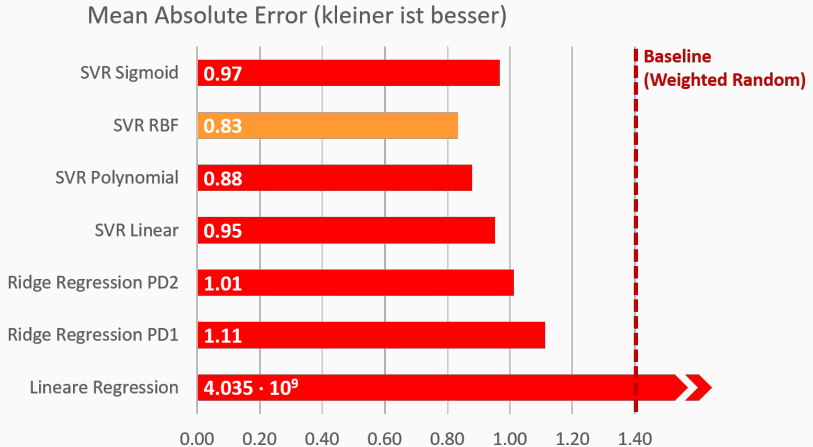
- Projekt: **Elasticsearch**
 - > 6 Jahre
 - ~ 22 000 Commits
 - ~ 10 000 Issues (GitHub)
- Trainingsset: 01.10.2014 - 31.12.2014
 - 620 Commits
 - 3 414 Dateiversionen
- Testset: 01.01.2015 - 31.01.2015
 - 165 Commits
 - 950 Dateiversionen

Workflow der ML-Pipeline



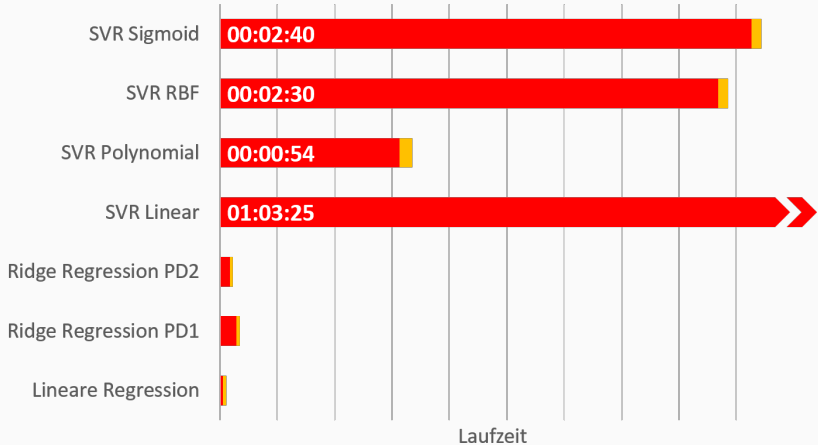
Resultate

Vergleich von ML-Algorithmen



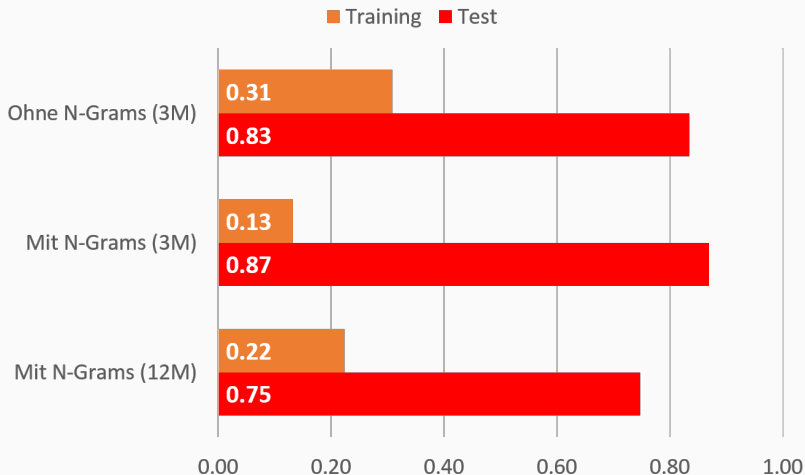
Laufzeit von ML-Algorithmen

Laufzeit pro ML-Konfiguration (kleiner ist besser)



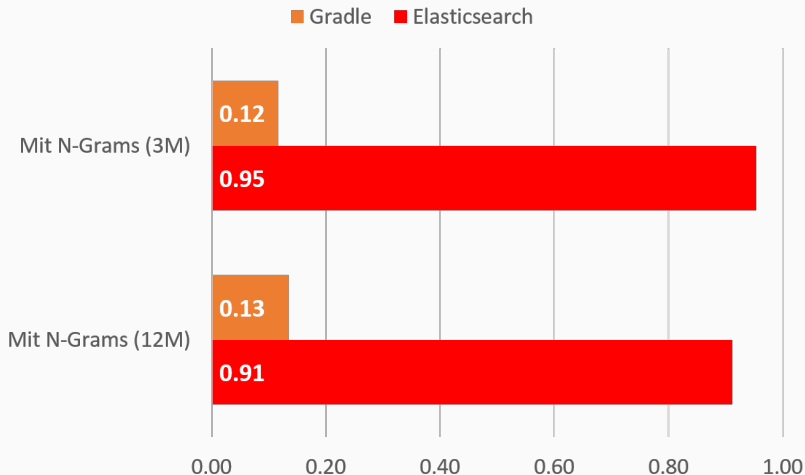
Performance von N-Grams

Median Absolute Error (kleiner ist besser)



Vergleich von Projekten

Median Absolute Error (kleiner ist besser)



Zusammenfassung

- Tool-Set für projektbezogenes Lernen
- Grundlage für weitere Arbeiten
- Textanalyse-Features
- einige Experimente

- Heatmap
- Fehlerlokalisierung
- Andere Programmiersprachen
- Cloud-Anwendung

Fragen und Diskussion
