



Fakultät für Informatik

Studiengang Software- und Systems-Engineering

# Erkennung von Design Patterns in Quellcode durch Machine Learning

Master Thesis

von

Mehmet Aslan

Datum der Abgabe: tt.mm.jjjj

Erstprüfer: Prof. Dr. Marcel Tilly

Zweitprüfer: Prof. Dr. Kai Höfig

#### EIGENSTÄNDIGKEITSERKLÄRUNG / DECLARATION OF ORIGINALITY

Hiermit bestätige ich, dass ich die vorliegende Arbeit selbständig verfasst und keine anderen als die angegebenen Hilfsmittel benutzt habe. Die Stellen der Arbeit, die dem Wortlaut oder dem Sinn nach anderen Werken (dazu zählen auch Internetquellen) entnommen sind, wurden unter Angabe der Quelle kenntlich gemacht.

*I declare that I have authored this thesis independently, that I have not used other than the declared sources / resources, and that I have explicitly marked all material which has been quoted either literally or by content from the used sources.*

Rosenheim, den tt.mm.jjjj

Vor- und Zuname

# Kurzfassung

text

Schlagworte:

# Inhaltsverzeichnis

<b>1</b>	<b>Motivation</b>	<b>1</b>
1.1	Einführung in Design Patterns . . . . .	2
1.2	Untersuchungsfragen . . . . .	2
<b>2</b>	<b>Literaturrecherche</b>	<b>4</b>
2.1	Design Patterns in der Software-Entwicklung . . . . .	5
2.1.1	Design Pattern Katalog . . . . .	6
2.1.2	Rollenkatalog . . . . .	7
2.2	Herausforderungen und Probleme bei der Erkennung von Design Patterns . . . . .	8
2.2.1	Variabilität der Implementierung . . . . .	8
2.2.2	Steigende Komplexität und Skalierbarkeit des Software-Systems . . . . .	8
2.2.3	Iterative Evolution des Quellcodes . . . . .	8
2.2.4	Mangel an expliziter Dokumentation . . . . .	8
2.3	Alternative Ansätze . . . . .	9
2.4	Angewendete Ansätze mit Maschine Learning . . . . .	10
2.5	Geeignete Datensätze . . . . .	11
2.5.1	Verfügbare gelabelte Datensätze . . . . .	11
2.5.2	Argumentieren des Datensatzes mit synthetischen Daten . . . . .	11
2.6	Feature Engineering . . . . .	12
2.6.1	Transformation von Quellcode in ASTs . . . . .	12
2.6.2	Extraktion von Software-Metriken . . . . .	12
2.6.3	Umwandlung von Quellcode in textuelle Form . . . . .	12
2.7	Betrachte Multiclass-Klassifizierer . . . . .	13
2.7.1	Ein-Modell Architekturen . . . . .	13
2.7.2	Mehr-Modell Architekturen . . . . .	13
2.8	Metriken für Klassifizierer . . . . .	14
2.8.1	F1 . . . . .	14
2.8.2	Recall . . . . .	14
2.8.3	Precision . . . . .	14
2.8.4	ROC . . . . .	14
2.9	Angewendete Technologien, Frameworks und Bibliotheken . . . . .	15
<b>3</b>	<b>Methodologie</b>	<b>16</b>
3.1	Verwendeter Datensätze . . . . .	17
3.2	Extrahierte Features . . . . .	18
3.3	Angewendete Multiclass-Klassifizierer . . . . .	19
3.4	Evaluation des trainierten Modells . . . . .	20
3.5	Klassifizierung . . . . .	21
3.5.1	Klassifizierung von Rollen in Design Patterns . . . . .	21
3.5.2	Klassifizierung von Design Patterns durch Rollen . . . . .	21
3.6	Evaluation der Methodologie . . . . .	22

<b>4</b>	<b>Zukünftige Aussichten</b>	<b>23</b>
<b>A</b>	<b>Erstes Kapitel des Anhangs</b>	<b>24</b>
	<b>Literaturverzeichnis</b>	<b>25</b>

# **Abbildungsverzeichnis**

# **Tabellenverzeichnis**

# **1 Motivation**



## 1.1 Einführung in Design Patterns

Entwurfsmuster oder auf Englisch ‘Design Patterns’ sind bewährte Lösungsansätze für wiederkehrende Probleme, die bei der Konzeption der Software-Architektur oder während der Implementierung der Software eingesetzt werden kann. Dabei dienen diese Entwurfsmuster als eine Art Blaupause, die es Software-Entwicklern ermöglicht, erprobte Lösungsstrategien für häufig auftretende Probleme in der Software-Entwicklung anzuwenden. Durch den Einsatz von etablierten Entwurfsmustern können Software-Entwickler für die Software bei korrekter Anwendung unter anderem erhöhte Wartbarkeit, Wiederverwendbarkeit von Komponenten, Verständlichkeit und Skalierbarkeit ermöglichen das wiederum in qualitativ besserer Software resultiert. Dabei sollte beachtet werden, dass Design Patterns als Vorlage zu betrachten sind. Je nach Einsatzgebiet muss die Anwendung des Entwurfsmusters evaluiert und für den konkreten Fall individualisiert werden. Deshalb existiert keine universelle anwendbare Iteration eines Design Patterns, die unabhängig von Anwendungskontext eingesetzt werden kann. Dies resultiert in variierenden Anwendung von Entwurfsmustern abhängig von jeweiligem Einsatzgebiet. Im weiteren Entwicklungszyklus der Software werden durch neue oder geänderte Anforderungen bereits eingesetzte Implementierungen von Entwurfsmustern modifiziert, entfernt oder neue werden hinzugefügt. Währenddessen besteht die Gelegenheit, dass durch mangelnder Dokumentation oder anderer Gründe die Entscheidungen, weshalb Entwurfsmuster so eingesetzt sind wie es eingesetzt worden, verloren gehen. Dadurch besteht die Gefahr, dass angewendete Design Patterns im weiteren Verlauf derer Entwicklung nicht mehr wiederzuerkennen sind. Aus diesem Grund ist die Etablierung eines Prozesses von Vorteil, das in der Lage ist, Implementierungen von Entwurfsmustern aus einem Software-System zu extrahieren und dieses konkret benennen. Vor allem der Einsatz von Maschine Learning für die Klassifizierung ist hier vorteilhaft, wodurch das Potenzial besteht, vorher nicht gesehene Implementierung von Design Patterns zu erkennen. Durch solch einen Prozess können durch die Erkennung von eingesetzten Entwurfsmustern auf konkrete und verlorenen gegangene Design-Entscheidungen zurückgeschlossen werden, welche zukünftige Design-Entscheidungen für das Software-System beeinflussen können. Der Fokus dieser Arbeit besteht daran, solch ein Prozess zu etablieren, welches für ein gegebenes Set von Quellcode-Dateien mithilfe von Maschine Learning einem potenziellen Entwurfsmuster zuzuteilen.

## 1.2 Untersuchungsfragen

Das Ziel dieser Arbeit besteht aus der Etablierung eines Prozesses, womit durch Einsatz von Maschine Learning für ein Set von Quellcode-Dateien ein Design Pattern zuzuordnen. Um solch ein Prozess zu entwickeln, werden in Kontext dieser Arbeit folgende Fragen beantwortet:

1. Welche Design Patterns werden berücksichtigt?
2. Was für ein Datensatz eignet sich für solch ein Prozess?
3. Wonach wird exakt klassifiziert?
4. Welche Merkmale, die aus Quellcode-Dateien extrahierbar sind, eignen sich für Klassifizierung Maschine Learning Modelle?

## *1 Motivation*

5. Welche Klassifizierer eignen sich?
6. Wie ist das Endresultat zu beurteilen?

## **2 Literaturrecherche**

## 2.1 Design Patterns in der Software-Entwicklung

Entwurfsmuster definieren gängige Lösungsblaupausen für häufig auftretende Probleme in Software-Entwicklung, vor allem in der Design-Phase der Architektur des Software-Systems als auch während der konkreten Implementierung. Jedoch sind diese als Schablone zu verstehen, die für den jeweiligen Einsatzfall angepasst werden müssen. Ein Werk, das das Verständnis von Design Patterns für das objektorientierte Programmieren maßgeblich geprägt ist, ist das von Gamma et al. verfasste Werk "Design Patterns: Elements of Reusable Object-Oriented Software". In diesem wird ein Katalog von 23 Entwurfsmustern definiert, welche in drei Kategorien aufgeteilt. Dieser Katalog wird von Software-Entwicklern als "Gang of Four" Entwurfsmuster bezeichnet. Gamma et al. definieren folgende Elemente für die Identifikation eines Entwurfsmusters: [Gam94, S. 3]

- **Pattern Name:** Der Name des Entwurfsmusters beschreibt in wenigen Worten, welches das zu lösende Problem, die Lösung und welche Folgen dessen Einsatz mit sich bringt. Durch die Einführung eines Bezeichners wird eine Schicht der Abstraktion hinzugefügt, welches das Verständnis und Dokumentation des Design Patterns vereinfacht.
- **Problem:** Das Problem beschreibt, wo das Entwurfsmuster angewendet werden soll. Dabei kann es sich um ein konkretes Entwurfsproblem, Klassen- oder Objektstrukturen oder eine Liste von Bedingungen darstellen, die zu erfüllen sind.
- **Solution:** Das Lösungselement beschreibt die Beziehungen, Verantwortlichkeiten und Zusammenarbeit der einzelnen Elemente, die die Struktur des Design Patterns definieren. Dabei werden diese Elemente in Objekte und Klassen, die die Grundbausteine der objektorientierten Programmierung repräsentieren, aufgeteilt und deren Interaktionen miteinander stellen die Verantwortlichkeiten und Beziehungen dar.
- **Consequences:** Die Folgen diskutieren, wie der Einsatz des betrachteten Entwurfsmusters sich auf das Software-System einwirkt und welche Vor- und Nachteile dadurch resultieren. Diese beeinflussen unter anderem die Zeit- und Speicherkomplexität, Erweiterbarkeit, Flexibilität und Portabilität des Software-Systems.

Im Kontext dieser Arbeit werden zu klassifizierende Strukturen, die potenziell einem Design Pattern zugeordnet werden können, als Mikroarchitekturen bezeichnet, die aus einer Menge von interagierenden Komponenten bestehen, denen je eine Rolle zugewiesen wird. Die jeweilige Rolle beschreibt, welche Funktionalität und Verantwortung diese im Kontext der Mikroarchitektur übernimmt und wie diese mit anderen Komponenten interagiert. Als Komponenten mit Rollen werden hier konkrete bzw. abstrakte Klassen oder Schnittstellen definiert, die die erforderliche Rolle im Rahmen der Mikroarchitektur erfüllen.

Im weiteren Verlauf dieser Sektion werden die drei erwähnten Entwurfsmusterkategorien erläutert und zu dem werden im Kontext dieser Arbeit betrachte spezifische Design Patterns genauer betrachtet.

### 2.1.1 Design Pattern Katalog

#### Creational Design Patterns

Die Kategorie der Creational Design Patterns oder Erzeugungsentwurfsmuster beschäftigt sich mit der Abstraktion des Prozesses der Initialisierung[Gam94, S. 81]. Entwurfsmuster dieser Kategorie fokussieren sich auf die Unabhängigkeit wie Objekte erstellt, zusammengesetzt und repräsentiert werden. Die Entwurfsmuster dieser Kategorie mit Fokus auf Klassen nutzen den Mechanismus der Vererbung, um zu beeinflussen, wie Komponenten instantiiert werden, während dahingegen Design Patterns mit einem Fokus auf Objekten die Instantiierung auf andere Objekte delegieren. Creational Design Patterns werden dann bedeutend, wenn mit steigender Komplexität des Software-Systems sich von Vererbung distanziert wird und die Komposition aus einzelnen definierten Objekt mehr an Bedeutung gewinnt[Gam94, S. 81]. Dabei wird das Verhalten einer Komponente auf eine Menge von einzelnen kleinere Objekten delegiert und durch Zusammensetzung innerhalb der Komponente und deren Interaktion das erwünschte Verhalten erzeugt. Dadurch wird die Instantiierung von Software-Komponenten komplexer, da die Instantiierung von mehreren Objekten koordiniert werden muss. Creational Design Patterns liefern hierbei Hilfestellung, weil die exakte Komposition der konkreten Objekte, die Teil der zu instantiierenden Komponente sind, und der exakte Prozess der Instantiierung im Inneren des Entwurfsmusters verborgen werden. Nach außen hin sind dahingegen nur die Schnittstellen sichtbar, die die Komponente zur Verfügung stellt, während dessen interne Logik die Ausführung auf andere Objekte delegiert. Im Kontext dieser Arbeit werden folgende Entwurfsmuster aus der Kategorie der Creational Design Patterns betrachtet:

#### Structural Design Patterns

Structural Design Patterns oder Strukturentwurfsmuster fokussieren sich darauf, wie einzelne Klassen und Objekte zusammengesetzt werden können, um größere Strukturen zu erzeugen[Gam94, S. 137]. Entwurfsmuster dieser Kategorie sind vorteilhaft, wenn unabhängig voneinander entwickelte Klassen oder Objekte aus verschiedenen Bibliotheken oder Frameworks miteinander interagieren müssen. Anstatt konkrete Implementierung und Schnittstellen zu nutzen, bedienen sich Structural Design Patterns der Komposition aus Objekten, um neue Funktionalitäten zur Verfügung zu stellen[Gam94, S. 137]. Die dadurch gewonnene Flexibilität ermöglicht das Ändern der Zusammensetzung des Objektes dynamisch zu der Laufzeit, welches mit statischer Komposition durch Klassen nicht möglich ist.[Gam94, S. 137]. Im Kontext dieser Arbeit werden folgende Entwurfsmuster aus der Kategorie der Structural Design Patterns betrachtet:

#### Behavioral Design Patterns

Behavioral Design Patterns oder Verhaltensentwurfsmuster konzentrieren sich auf Algorithmen und der Zuweisung von Verantwortlichkeiten zwischen Objekten[Gam94, S. 221]. Dabei wird nicht nur Struktur der Entwurfsmuster betrachtet, sondern auch die Kommunikation und Interaktion der Objekte, die Teil des Entwurfsmusters sind. Charakteristisch für Design Patterns dieser Kategorie ist der Fokus auf Verknüpfung der einzelnen Teilobjekte des Entwurfsmusters, anstatt des Kontrollflusses, welcher zur Laufzeit schwer nachvollziehbar sein kann[Gam94, S. 221]. Im Kontext dieser Arbeit werden folgende Entwurfsmuster aus der Kategorie der Behavioral Design Patterns betrachtet:

### **2.1.2 Rollenkatalog**

## **2.2 Herausforderungen und Probleme bei der Erkennung von Design Patterns**

### **2.2.1 Variabilität der Implementierung**

Design Patterns stellen in der Software-Entwicklung bewährte Lösungsmuster für bereits begegnete Herausforderung dar. Aufgrund der abstrakten und wiederverwendbaren Natur der Entwurfsmuster, muss für diese eine konkrete Implementierung definiert werden, die von dem Einsatzfall, Kontext und anderen Faktoren wie verwendeter Programmiersprache, Bibliotheken und Erfahrungsstand des Software-Entwicklers. Dadruch, dass jedes Entwurfsmuster einen konzeptionellen Rahmen darstellt und jede Implementierung von nicht statischer Außenfaktoren beeinflusst wird, resultiert dies in einem breiten Spektrum an Implementierungen für ein gegebenes Entwurfsmuster. Aus diesem Grund ist eine Definition einer starren Definition eines Design Patterns, was als Startpunkt und Referenz für die Erkennung des jeweiligen Entwurfsmusters dienen könnte, nicht möglich. Deshalb ist eine definitive Antwort auf die Frage, ob eine betrachte Mikroarchitektur eine Instanz eines Entwurfsmusters, nicht beantwortbar, weshalb die Antwort von automatisierten Prozessen von Design Patterns eher mit einem Wert besteht, welches die Ähnlichkeit zu einem Design Pattern beschreibt.

### **2.2.2 Steigende Komplexität und Skalierbarkeit des Software-Systems**

### **2.2.3 Iterative Evolution des Quellcodes**

### **2.2.4 Mangel an expliziter Dokumentation**

### **2.3 Alternative Ansätze**



## **2.4 Angewendete Ansätze mit Maschine Learning**

## **2.5 Geeignetes Datensätze**

### **2.5.1 Verfügbare gelabelte Datensätze**

### **2.5.2 Argumentieren des Datensatzes mit synthetischen Daten**

## **2.6 Feature Engineering**

### **2.6.1 Transformation von Quellcode in ASTs**

### **2.6.2 Extraktion von Software-Metriken**

### **2.6.3 Umwandlung von Quellcode in textuelle Form**

## **2.7 Betrachte Multiclass-Klassifizierer**

### **2.7.1 Ein-Modell Architekturen**

### **2.7.2 Mehr-Modell Architekturen**

## **2.8 Metriken für Klassifizierer**

### **2.8.1 F1**

### **2.8.2 Recall**

### **2.8.3 Precision**

### **2.8.4 ROC**

## **2.9 Angewendete Technologien, Frameworks und Bibliotheken**

## **3 Methodologie**

### **3.1 Verwendeter Datensätze**



## **3.2 Extrahierte Features**

### **3.3 Angewendete Multiclass-Klassifizierer**

### **3.4 Evaluation des trainierten Models**

## **3.5 Klassifizierung**

### **3.5.1 Klassifizierung von Rollen in Design Patterns**

### **3.5.2 Klassifizierung von Design Patterns durch Rollen**

### **3.6 Evaluation der Methodologie**

## **4 Zukünftige Aussichten**

## **A Erstes Kapitel des Anhangs**

Wenn Sie keinen Anhang benötigen, dann bitte einfach rausnehmen.

# Literaturverzeichnis

- [Gam94] E. Gamma, R. Helm, R. Johnson und J. M. Vlissides. *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley Professional, 1 Aufl., 1994.